

A Critical Review of Large Language Model on Software Engineering: An Example from ChatGPT and Automated Program Repair

Quanjun Zhang
quanjun.zhang@smail.nju.edu.cn
State Key Laboratory for Novel
Software Technology, Nanjing
University
Nanjing, Jiangsu, China

Chunrong Fang*
fangchunrong@nju.edu.cn
State Key Laboratory for Novel
Software Technology, Nanjing
University
Nanjing, Jiangsu, China

Tongke Zhang
201250032@smail.nju.edu.cn
State Key Laboratory for Novel
Software Technology, Nanjing
University
Nanjing, Jiangsu, China

Bowen Yu
201250070@smail.nju.edu.cn
State Key Laboratory for Novel
Software Technology, Nanjing
University
Nanjing, Jiangsu, China

Zhenyu Chen
zychen@nju.edu.cn
State Key Laboratory for Novel
Software Technology, Nanjing
University
Nanjing, Jiangsu, China

Juan Zhai
juanzhai@umass.edu
Manning College of Information &
Computer Sciences, University of
Massachusetts
Amherst, MA, United States

Weisong Sun
weisongsun@smail.nju.edu.cn
State Key Laboratory for Novel
Software Technology, Nanjing
University
Nanjing, Jiangsu, China

ABSTRACT

Large Language Models (LLMs) have been gaining increasing attention and demonstrated promising performance across a variety of Software Engineering (SE) tasks, such as Automated Program Repair (APR), code summarization, and code completion. For example, ChatGPT, the latest black-box LLM, has been investigated by numerous recent research studies and has shown impressive performance in various tasks. However, there exists a potential risk of data leakage since these LLMs are usually close-sourced with unknown specific training details, *e.g.*, pre-training datasets.

In this paper, we seek to review the bug-fixing capabilities of ChatGPT on a clean APR benchmark with different research objectives. We first introduce EVALGPTFIX, a new benchmark with buggy and the corresponding fixed programs from competitive programming problems starting from 2023, after the training cutoff point of ChatGPT. The results on EVALGPTFIX show that ChatGPT is able to fix 109 out of 151 buggy programs using the basic prompt within 35 independent rounds, outperforming state-of-the-art LLMs

CodeT5 and PLBART by 27.5% and 62.4% prediction accuracy. We also investigate the impact of three types of prompts, *i.e.*, problem description, error feedback, and bug localization, leading to additional 34 fixed bugs. Besides, we provide additional discussion from the interactive nature of ChatGPT to illustrate the capacity of a dialog-based repair workflow with 9 additional fixed bugs. Overall, our experiments demonstrate that ChatGPT is able to fix a total of 143 bugs in EVALGPTFIX, indicating the potential of ChatGPT in repairing real-world buggy programs. Inspired by the findings, we further pinpoint various challenges and opportunities for advanced SE study equipped with such LLMs (*e.g.*, ChatGPT) in the near future. More importantly, our work calls for more research on the reevaluation of the achievements obtained by existing black-box LLMs across various SE tasks, not limited to ChatGPT on APR.

CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

KEYWORDS

Automated Program Repair, Large Language Model, AI4SE

ACM Reference Format:

Quanjun Zhang, Tongke Zhang, Juan Zhai, Chunrong Fang, Bowen Yu, Weisong Sun, and Zhenyu Chen. 2023. A Critical Review of Large Language Model on Software Engineering: An Example from ChatGPT and Automated Program Repair. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/nmnnnnn.nmnnnnn>

* Chunrong Fang is the corresponding author.

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 ACM 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](https://permissions.acm.org).

Conference acronym 'XX, June 03–05, 2023, Woodstock, NY

© 2023 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/23/06...\$15.00
<https://doi.org/10.1145/nmnnnnn.nmnnnnn>

1 INTRODUCTION

The scale of modern software systems has been continuously expanding in recent years, leading to a significant surge in the number of bugs within these systems [15, 39]. Manual debugging is one of the critical software development activities, which usually requires a substantial investment of time and human resources to keep the software well-maintained [4, 5]. In order to reduce the costs of repairing such detected bugs, Automated Program Repair (APR) is proposed [27], aiming at generating correct patches automatically to minimize human involvement in the manual debugging process.

In the literature, existing APR techniques can be categorized into traditional and learning-based ones [14, 40]. Among traditional APR techniques, template-based APR, which mainly relies on pre-defined repair templates to transform buggy code into the correct one, has been considered as state-of-the-art [52, 53]. However, it is challenging to fix unseen bugs that fall outside the scope of pre-defined templates. On the other hand, learning-based APR is able to learn the bug-fixing patterns automatically from a large code repository on top of the advance of deep learning [8, 9]. Learning-based APR usually leverages Neural Machine Translation (NMT) model to translate a code sequence from a source language (i.e., buggy code snippets) into a target language (i.e., correct code snippets) [47]. Despite addressing the limitations of template-based APR and demonstrating promising results, the performance of learning-based APR relies on the quality and quantity of the training data [53].

More recently, Large Language Models (LLMs) are gaining increasing attention due to their powerful programming language processing capabilities in various Software Engineering (SE) tasks [13, 37, 52]. These LLMs are usually trained with a pre-training-and-fine-tuning mechanism [12, 49], i.e., pre-trained by self-supervised training on a large-scale unlabeled corpus to derive generic knowledge, and fine-tuned by supervised training on a limited labeled corpus to benefit a specific downstream task. Among existing LLMs, ChatGPT [43] is widely regarded as one of the most popular language models today and is being studied by researchers from numerous domains, such as code summarization [46], code generation [30] and test generation [57]. In particular, ChatGPT is a prompt-based LLM equipped with Reinforcement Learning from Human Feedback that can interact with users through human-like dialogues. In the domain of APR, researchers have attempted to directly utilize ChatGPT in generating correct patches and have yielded promising results [45, 54]. For example, Sobania *et al.* [45] evaluate the bug-fixing capabilities of ChatGPT and find ChatGPT is able to fix 31 out of 40 bugs on QuixBugs benchmark. Xia *et al.* [54] employ ChatGPT in a conversational manner to fix 114 and 48 bugs on the Defects4J-v1.2 and Defects4J-v2.0 benchmark, and all the 40 bugs from QuixBugs benchmark, significantly outperforming state-of-the-art APR techniques.

In spite of the remarkable performance, there are some concerns with the well-known dataset used to evaluate ChatGPT for APR. ChatGPT is trained on vast amounts of data from the internet, which may contain data in the commonly-chosen datasets for APR (e.g., Defects4J [23] and QuixBugs [29]). It is difficult to ensure whether or not the evaluation dataset has not been seen by ChatGPT during training. For example, when we ask ChatGPT whether it is

aware of Defects4J, as shown in Figure 1, ChatGPT gives an affirmative answer and can further list the projects present in Defects4J. If ChatGPT has previous knowledge of the dataset, employing it to fix bugs from the dataset might not well reflect its fixing ability since it may already be aware of the bug-fixing patches. The concern exists in other LLMs and code-related tasks as well. A similar example shows¹ that GPT-4 is able to solve all 10 code competition problems until 2021, which is the training cutoff of the model, while none is solved correctly for 10 problems after that date instead. We also find that ChatGPT can directly provide complete descriptions and the corresponding solution by simply providing it with the number of a programming problem in LeetCode (presented in our online repository [2]). Similarly, Karmakar *et al.* [24] highlight the memorization issue of Codex, showing its ability to generate accurate code outputs only with the first sentence of the problem description as a prompt, even in the absence of clear task objectives. Considering the fact there exist a quite number of black-box LLMs for which no architecture or training data information has been released. The data leaking on such LLMs is a significant concern when it comes to evaluating their performance in some code-related tasks in the SE community, such as APR [52] and assertion generation[41].

This paper. In our work, we attempt to raise the important concern about data leakage, which has been an overlooked issue in the SE community, when such black-box LLMs are applied to some code-related tasks. We select ChatGPT and APR as representative examples of LLMs and SE tasks, respectively. In particular, we construct a new dataset EVALGPTFIX with buggy and correct code nippets from a competitive programming website Atcoder. We crawl users' submissions for competitions during 2023. As ChatGPT states its knowledge cutoff is in September 2021, thus we confirm it does not have access to samples in EVALGPTFIX. We then design the following three research questions to evaluate ChatGPT's bug-fixing ability in our experiments:

(RQ1) The effectiveness of ChatGPT on EVALGPTFIX.

Results: In this RQ, we aim to investigate how ChatGPT behaves in repairing bugs when we present it with the buggy code. Although ChatGPT could learn more from detailed prompts, we only give a prompt containing the bug and an instruction that simply asks to fix the bug. This is similar to the input for common APR methods, so we can compare the repair effectiveness of ChatGPT with other tools. We require ChatGPT to try to fix 151 bugs chosen from two latest AtCoder competitions in EVALGPTFIX, and count the number of bugs that ChatGPT is able to fix to explore to what extent can ChatGPT generate correct patches.

(RQ2) What is the effect of different prompts on the repair performance of ChatGPT?

Results: ChatGPT is a prompt-based language model, whose response is largely dependent on how it is prompted. Through modifying prompts, we can give more information relating to the bugs. After gaining such information, ChatGPT has the potential to fix more bugs. To find out what prompt helps in bug fixing, we design three advanced prompts which separately include the programming problems that the programs target at, the exact lines where the bugs are in, and what

¹<https://github.com/hitz-zentroa/lm-contamination>

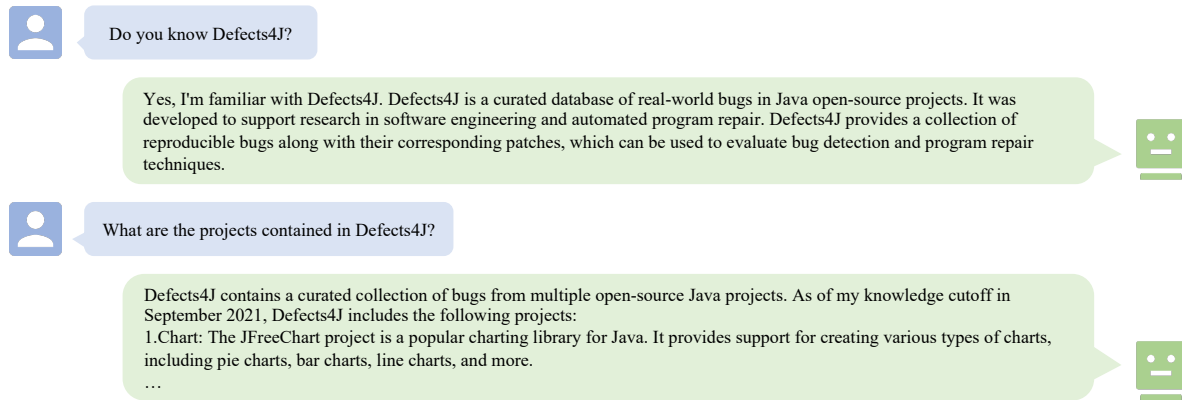


Figure 1: Dialogue with ChatGPT about its knowledge of Defects4J Benchmark

type of error the bugs have. We apply these prompts to the bugs not fixed in RQ1 and experiment on how many more bugs can be fixed by each prompt.

(RQ3) Can dialogues help ChatGPT in improving repair performance?

Results: ChatGPT is created to engage in dialogues with users, and while answering a query, it has a memory of previous conversations, from which it can adjust its responses. For the bugs that ChatGPT fails to fix when only a single round of dialogue is performed, we continue the dialogues by informing ChatGPT of what is wrong with the patches it has generated. The intuition is that ChatGPT can learn from the implausible fixes and gradually turn its attention to the correct patch with the accumulative dialogues.

Overall, our work confirms that ChatGPT has an excellent performance in fixing bugs from the dataset EVALGPTFIX. In the case of providing the basic prompt, ChatGPT fixes 109 bugs out of the 151 bugs. When programming problem descriptions, error messages and bug locations are added to the prompt, respectively 18, 25 and 10 more bugs are fixed. Moreover, by conducting dialogues, ChatGPT fixes 9 bugs that are neither fixed with the basic prompt nor with the prompt including error information. The results show that the repair ability of ChatGPT can benefit from prompts that are more in detail and deeper dialogues.

Novelty & Contributions. To sum up, the main contributions of this paper are as follows:

- **Overlooked Issue.** We reveal an important concern when evaluating the recent ChatGPT in repairing software bugs with commonly-adopted benchmarks, i.e., *the data leakage issue*. More importantly, the issue potentially exists in a broader range of other code-related tasks and black-box LLMs, and has been consistently overlooked by the SE community. Thus, the issue may lead to a significant bias in previous research works that employ black-box LLMs (such as ChatGPT and Codex) without assessing any training details, such as pre-training datasets and model architectures.

- **Clean Benchmark.** We construct a new APR benchmark EVALGPTFIX from a competitive programming website At-coder. EVALGPTFIX contains 151 pairs of bugs and fixes in Java, which come from failing and accepted programming submissions in 2023 to ensure that ChatGPT has not seen the specific code snippets presented in this dataset.
- **Extensive Study.** We conduct an in-depth empirical study of how ChatGPT are applied to automated program repair. Specifically, our study is three-fold: (1) a systematic comparison between ChatGPT and state-of-the-art LLMs, indicating that ChatGPT can outperform CodeT5 and PLBART; (2) an extensive evaluation to analyze the impact of different prompts; (3) an additional discussion about the impact of dialogue-based repair workflow for ChatGPT.
- **Challenge and Opportunity** We discuss current pressing challenges and forward-looking directions on applying ChatGPT and more advanced black-box LLMs for future program repair and other SE studies.

Open Science. To support the open science community, we release the studied dataset, scripts (i.e., data processing, model training, and model evaluation), and related models in our experiment for replication and future research [2].

2 BACKGROUND

2.1 Automated Program Repair

Automated Program Repair (APR) is raised to generate candidate patches automatically for the buggy code snippets, so as to reduce the time and cost of manual debugging [55, 60]. There are mainly two types of APR techniques, i.e., traditional and learning-based ones.

Traditional APR can be classified into three categories: heuristic-based [27, 35, 56], constraint-based [11, 36, 38], and template-based [26, 31, 32]. Among them, template-based APR has shown promising performance in the bug-fixing task. Template-based APR utilizes pre-defined fix templates, which are patterns of code changes commonly applied in debugging activities, to generate possible patches for specific bugs. Despite of its significant ability in program repair, template-based APR has limitations in both fix templates and donor

code. It cannot fix bugs requiring action beyond the fix templates, and with a proper fix template, some bugs still cannot be fixed because of a limited source of donor code [32].

The problems above can be solved by learning-based APR, which is a rising field that APR researchers are focusing on. APR problems are regarded as Neural Machine Translation (NMT) tasks that transform a piece of buggy code into the correct one. Learning-based APR leverages DL techniques to gain more insight into program repair behaviors from large code corpora. However, the effectiveness of learning-based APR can be easily affected by the quality of training data, which may contain code changes irrelevant to bug fixing and thus limit the performance of APR.

2.2 Large Language Model

Large Language Models (LLMs) are language models consisting of billions of parameters trained from a significant amount of data and have impressive performance in language processing tasks, including both natural languages and programming languages [12, 17, 18, 24, 49]. LLMs are typically based on Transformer architecture, where an encoder takes a variable-length input and turns it into a fixed-length vector, while a decoder transforms the encoded representation into a sequence of output. Based on the employed components, LLMs can be categorized into encoder-only, decoder-only, and encoder-decoder ones. Encoder-only models (e.g., BERT[10]) learn data representation through training objectives like Masked Language Modeling (MLM). Decoder-only models (e.g., GPT[6]) are trained to generate predictions for the next token given all the previous tokens. Encoder-decoder models (e.g., CodeT5 [49]) combine the encoder and the decoder, and are trained with the objective of recovering the corrupted input.

3 STUDY DESIGN

3.1 Research Questions

In this work, we explore the following research questions.

- RQ1:** What is the performance of ChatGPT in repairing buggy programs from EVALGPTFIX?
- RQ2:** How do the different prompts with additional program information affect the performance of ChatGPT?
- RQ3:** How does the interaction with dynamic execution feedback affect the performance of ChatGPT?

3.2 EVALGPTFIX Construction

There is evidence showing that ChatGPT already has knowledge about popular datasets (e.g., Defects4J and Quixbugs) widely used by current APR techniques. As a result, it seems not rigorous to research on ChatGPT's program repair ability based on these datasets as ChatGPT may find a correct patch for the given bug from its training data instead of fixing the bug by itself. To solve this problem, we construct a new dataset whose data is invisible to ChatGPT. We gain our data from AtCoder, a platform for programming contests. We extract the programming problems from contests in 2023 and get users' submissions of these problems as the source data. As ChatGPT is trained on data before September 2021, it has limited knowledge about our dataset, so applying ChatGPT to fix bugs from

EVALGPTFIX can better reflect ChatGPT's performance in the APR task.

① **Raw Data Collection.** We first crawl all the Java submissions of AtCoder programming contests starting from 2023. We focus on Java languages as it is the most targeted language in the APR community. The online judge results of the submissions can be divided into six types: 1) Accepted (AC); 2) Wrong Answer (WA); 3) Time Limit Exceeded (TLE); 4) Compilation Error (CE); 5) Runtime Error (RE); and 6) Memory Limit Exceeded (MLE).

② **Bug-Fixing Pairs Construction.** For all the submissions of a user on a single problem, we take the unaccepted submissions as the buggy program and the accepted submission as a corresponding correct program. Then we calculate the token difference between every pair of the buggy and correct programs, and only keep those with a difference of less than 6 tokens following existing the study [19]. This difference is calculated by tokenizing each program into a sequence of tokens and counting the number of token-level differences, including replacements, deletions, and insertions, between the two programs. This setting is based on the *competent programmer hypothesis*, i.e., seasoned programmers have the competence to produce programs that are nearly error-free, and that the majority of bugs can be rectified through minor modifications.

③ **Test Case Mining.** For each problem, the problem description in the HTML website includes a handful of illustrative input-output pairs that serve as examples. However, these test cases are not sufficient to validate the correctness of generated patches due to the overfitting problem in APR [59]. We further download all possible public test cases of all the problems in our dataset from the dedicated database² that posts all test cases of AtCoder problems. These test cases are manually created by domain experts, such as the programming problem setters, and serve as the test oracle in the backend of the website to assess the functional correctness of programs submitted by users.

④ **Static-based Filtering.** Furthermore, we remove repeating submissions as well as submissions with more than 500 code tokens, considering the limitation of repair models' ability to handle long code snippets (e.g., Tufano *et al.* [48] limit the maximum length of the buggy code to 100 tokens and CIRCLE [50] truncates the first 512 tokens of the buggy code). We also delete the comments in the code, as comments can provide information about the function of the buggy code, and can affect the judgment of APR tools' capability in fixing code without other hints. We also observe that some programs have a compilation error only because the class name is not written as "Main". This is irrelevant to the logic of the code itself, so such data is deleted from our dataset.

⑤ **Dynamic-based Filtering.** We execute all remaining submissions against every test case associated with the respective problem. The time limit and a memory limit of running a test case are respectively set as 10 seconds and 1MB following previous work [44]. In a pair (s_1, s_2) where s_1 represents the buggy code and s_2 represents the fixed one, if any of the following three conditions happen, the pair will be removed: (1) s_1 passes all the test cases of its corresponding problem; (2) The bug type of s_1 does not match the one given on AtCoder website (e.g., s_1 is found to produce a "Wrong

²https://www.dropbox.com/sh/arnpe0ef5wds8cv/AAAk_SECQ2Nc6SVGii3rHX6Fa?dl=0

Answer" error by running test cases but it gets a label of "Compilation Error" on AtCoder); and (3) s2 fails any of the test cases of its corresponding problem. Besides, due to differences in our local device and AtCoder platform environments, we exclude bugs with a type of MLE, resulting in four types of bugs in EVALGPTFIX (*i.e.*, WA, TLE, CE, and RE).

⑤ **Benchmark Statistics.** After all pre-processing phrases, we get 151 pairs of bug-fixing of Java programs for 15 programming problems from the two latest programming contests when we conduct the work, *i.e.*, Beginner Contest 297 and 298.

3.3 ChatGPT Setup

We conduct our experiment based on the API of ChatGPT with the model gpt-3.5-turbo released by OpenAI. Considering that ChatGPT will generate different responses when it is queried by the same input several times, with every prompt we send the request to ChatGPT repeatedly to improve the possibility of generating more correct patches.

3.4 Evaluation Metrics

We evaluate whether the patch generated by ChatGPT is correct by running it against the test suite. All the test cases are downloaded from the website of AtCoder, and they are used to judge users' submissions during coding contests. Averagely, every programming problem has 38 test cases, which is enough to tell whether the program can correctly solve the problem. We run every candidate fix on the test cases of the problem, and if it passes all the test cases, it is regarded as a correct fix.

For every bug with a prompt, we ask ChatGPT for a fix continuously in a loop, and if in any of the rounds, ChatGPT can generate a correct patch for a bug, we consider ChatGPT to be able to fix the bug successfully, and the loop will be exited.

3.5 Compared Techniques

We consider the following two state-of-the-art LLMs as the baseline techniques.

- **CodeT5.** Wang *et al.* [49] introduce a pre-trained language model (*i.e.*, CodeT5) on top of the T5 architecture by incorporating the token type information. CodeT5 considers both unimodal (code only) and bimodal (code-text pairs) data for four pre-training tasks, *i.e.*, masked span prediction, masked identifier prediction, identifier tagging, and bimodal dual generation.

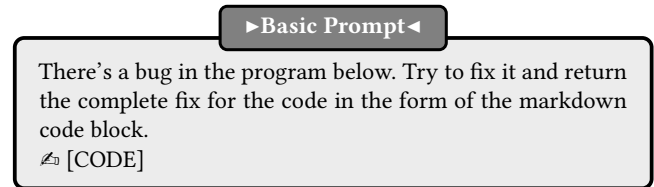
- **PLBART.** Ahmad *et al.* [1] introduce a pre-trained encoder-decoder model (*i.e.*, PLBART) on top of the BERT architecture to perform both program and language understanding and generation tasks. PLBART considers three denoising auto-encoding strategies to reconstruct an input text that is corrupted by a noise function in pre-training, *i.e.*, token masking, token deletion, and token infilling.

4 RESULTS AND ANALYSIS

4.1 RQ1: Effectiveness of ChatGPT

Design. In this RQ, we explore the repair ability of ChatGPT by presenting it with buggy programs and asking it to repair them. We only give the buggy code without any other information about the bugs to find out to what extent can ChatGPT localize and repair

bugs if no extra prompts are provided. The basic prompt designed for RQ1 is presented as follows, where [CODE] represents the buggy program to be fixed.



ChatGPT generates different responses even if it is prompted by the same sentences, so if it is not able to return the correct patch for a bug, there is a possibility that it will fix the bug when queried again. Therefore, asking ChatGPT to fix a bug only once cannot reflect its actual repair capability well. To solve this problem, for every bug, we repetitively send the same request to ChatGPT. If a patch for a bug can pass all the test cases, we stop asking ChatGPT to fix the bug again. In every round of query, we check whether any new bugs are fixed compared to the last round, and if no more bugs are fixed for three consecutive rounds, the process is stopped.

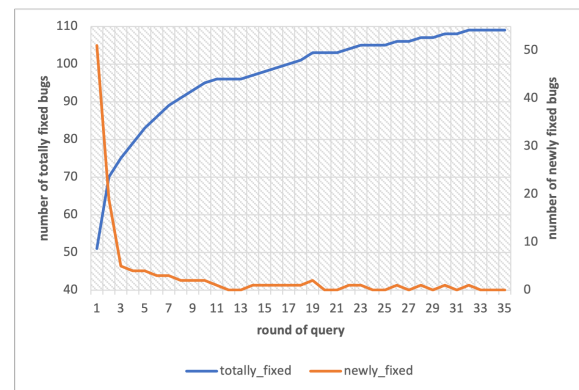


Figure 2: The number of totally fixed and newly fixed bugs in every round of query to ChatGPT

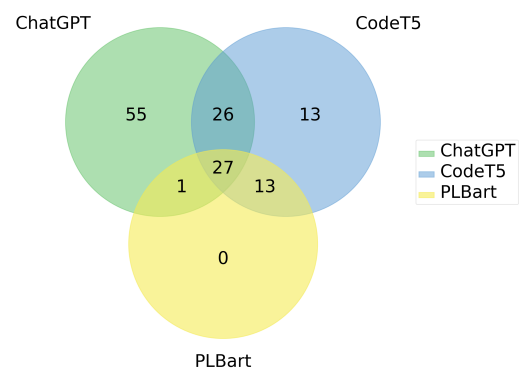


Figure 3: The overlap of bugs fixed by three models

```

1  ...
2  String S = sc.next();
3  boolean f1 = false;
4  boolean f2 = true;
5  for(int i=0; i<N; i++){
6  -   if(S[i]=='o'){
7  +   if(S.charAt(i)=='o'){
8     f1 = true;
9   }
10 -   if(S[i]=='x'){
11 +   if(S.charAt(i)=='x'){
12     f2 = false;
13   }
14 }
15 ...

```

Listing 1: An example of bug with syntax error fixed by basic prompt

Results. Figure 2 shows the number of bugs that ChatGPT is able to fix in each round of requests. We finally query ChatGPT for 35 independent rounds, in which the 33-35 rounds do not see any newly-fixed bugs. Among the 151 bugs from the two latest programming contests at AtCoder, 109 are successfully fixed. Generally, there is a decline in additional fixes with the increase of round number. In the first three rounds, respectively 51, 19, and 5 additional correct patches are generated, which is a sharp decrease. After that, the number begins to drop slowly, and vacillates between 0 to 2 since the 11th round. The decreasing trend only stopped after 35 rounds, indicating the randomness nature of ChatGPT. Therefore, we recommend future work to explore the issue of randomness in ChatGPT, which has been ignored in most previous studies, e.g., Sobania *et al.* [45] only repeats four times.

We then investigate the types of bugs fixed by ChatGPT. Among the 151 bugs in EVALGPTFIX, there exist 23, 4, 22, and 102 bugs of the four types CE, TLE, RE, and WA. ChatGPT is able to fix 22, 4, 11, and 72 of them, with the fixing percentage of 96%, 100%, 50%, and 71%. We find ChatGPT shows an impressive performance in localizing and fixing bugs with Compilation Error(CE), which is usually caused by syntax errors. Such errors are easy to identify as long as ChatGPT has knowledge of Java syntax, so it can fix the bug without considering the operating logic of the program, which can be much more intricate. For example, in the following bug shown in Listing 1, the array operator "`[]`" is used on a string, which is not allowed in Java. ChatGPT recognizes the error and replaces "`S[i]`" with "`S.charAt(i)`", and gives a correct patch.

Apart from simple syntax errors, ChatGPT can also identify and repair some logic errors in the programs. Listing 2 presents an example with a logic error that can be successfully fixed by ChatGPT. In the buggy program, the conditional expression "`i<t.length`" in the while statement is problematic, as in the next statement it can trigger `ArrayIndexOutOfBoundsException`, which is a kind of Runtime Error. ChatGPT changes the buggy expression into "`i<t.length-1`" to prevent the exception, so as to solve the problem in the code snippet.

We further compare the repair performance of ChatGPT with the other two state-of-the-art pre-trained models, i.e., CodeT5 and PLBART. The models are fed with the buggy code and then generate possible patches which are later run on the test cases. We fine-tune the selected models with the FixEval dataset [19], which contains 156k buggy and correct code submissions to competitive

```

1  ...
2  int[] t = new int[n];
3  ...
4  - while(i<t.length)
5  + while(i<t.length-1)
6  {
7    {
8      if(t[i+1]-t[i]<=d)
9    }
10 ...

```

Listing 2: An example of bug with logic error fixed by basic prompt

programming problems before 2021. The beam size is set as 50 for both models, which means 50 patches with highest possibility are produced for each bug. We find that CodeT5 fixes 79 bugs and PLBART only fixes 41 bug, 27.5% and 62.4% less than what is achieved by ChatGPT. Figure 3 presents the overlap of bugs fixed by the three models, showing that ChatGPT significantly fixes more bugs (i.e., 55 unique bugs) that the other two models are not able to fix (i.e., 13 unique bugs for CodeT5 and none for PLBART). The results indicate that ChatGPT notably outperforms existing LLMs with regard to repairing programming problems.

Answer to RQ1: The performance of ChatGPT in EVALGPT-Fix shows that: (1) there exists a significant randomness issue in ChatGPT, e.g., 35 independent rounds are required to achieve stable results; (2) ChatGPT is effective in fixing different types of bugs, e.g., 96%, 100%, 50% and 71% of CE, TLE, RE and WA bugs are correctly fixed; (3) ChatGPT is able to fix 109 bugs in EVALGPTFIX with a recall of 72.19%, 30 and 68 more than CodeT5 and PLBART.

4.2 RQ2: The Impact of Prompt

Design. We further add more details about the bug to the prompt given to ChatGPT, expecting that ChatGPT will gain more useful information from the prompts so that it can fix more bugs. We ask ChatGPT to fix the bugs that are not fixed when only the basic prompt is offered, as described in RQ1. We consider three types of additional bug information, including problem descriptions, error information, and bug locations. For the three types of bug information, we respectively design three kinds of prompts based on the original prompt, discussed as follows.

- **Problem descriptions** indicate what the programming problem aims to solve. The prompt with a problem is shown below, where [CODE] represents the buggy program and [PROBLEM] represents the coding problem that the code is submitted to. All the problem descriptions are obtained from the website of AtCoder, and consist of the background of the problem, the input to the program, and what the output is supposed to be like.

► Problem Description Prompt ◀

📌 There's a bug in the program below. Try to fix it and return the complete fix for the code in the form of the markdown code block.

📌 [CODE]

The program is to solve this problem:

📌 [PROBLEM]

• **Error information** informs ChatGPT of the type of error, the input that triggers the error, the expected correct output, and the actual wrong output of the buggy program. For the four types of bugs (i.e. CE, TLE, RE, and WA), the prompts are slightly different from each other, displayed as follows.

❶ **Compilation Error (CE)** is not triggered by any input as it happens while compiling, so in the case of a CE bug, we only tell ChatGPT that there is a compilation error in the code.

► **Compilation Error Prompt** ◀

There's a bug in the program below. Try to fix it and return the complete fix for the code in the form of the markdown code block.

📄 [CODE]

There's a Compilation Error in the code.

❷ For a **Time Limit Exceeded (TLE)** or **Runtime Error (RE)**, the input that causes the error as well as the expected output are added to the prompt. The prompt is described as follows, where [INPUT] represents the input that causes the test case to fail and [EXPECT] represents the output that should be printed by a correct program.

► **Exceeded/ Runtime Error Prompt** ◀

There's a bug in the program below. Try to fix it and return the complete fix for the code in the form of the markdown code block.

📄 [CODE]

The following input triggers a Time Limit Exceeded/ Runtime Error:

📄 [INPUT]

The expected output is:

📄 [EXPECT]

❸ For a bug with a **Wrong Answer (WA)** error, apart from the input and expected output, the prompt also contains the actual output of the buggy program. The prompt is described as follows, where [OUTPUT] represents the output of the buggy program given the input filled in "[INPUT]".

► **Wrong Answer Error Prompt** ◀

There's a bug in the program below. Try to fix it and return the complete fix for the code in the form of the markdown code block.

📄 [CODE]

The following input triggers a Wrong Answer error:

📄 [INPUT]

The expected output is:

📄 [EXPECT]

The actual output is:

📄 [OUTPUT]

• **Bug localization** show the suspicious lines where bugs are localized in [51]. We mark the buggy lines with a comment "//bug" and then ask ChatGPT to repair the code with a bug position. To obtain the buggy line, we compare the lines of every pair of bug

and fix. In the bug, the first line that is different from the line in the fix is considered as the buggy line. We add another short prompt just following the basic prompt to tell ChatGPT what the comment means. The prompt is designed as follows, where [CODE] represents the buggy code whose buggy line is followed by the comment "//bug".

► **Bug Localization Prompt** ◀

There's a bug in the program below. Try to fix it and return the complete fix for the code in the form of the markdown code block. The location of the bug is in or near the line with a comment "//bug".

📄 [CODE]

Results. In the 42 bugs that are not fixed in RQ1, 25, 18, and 10 more bugs are fixed when we separately add error information, problem description, and bug location to the basic prompt. This shows the promoting effect that more concrete prompts have on the repair performance of ChatGPT. Like in RQ1, we query ChatGPT several times until no bugs are newly fixed in continuously three rounds. Figure 4 shows how many additional bugs are fixed in every round. It only takes 12 and 13 rounds to prompt with problems and bug locations, but 27 rounds are executed when it comes to error information, notably more than the other two prompts. The possible reason is that the former two prompts help ChatGPT to focus on a smaller range of code, so the patches it generates in every round are relevantly stabler. By contrast, there is vacillation in the patches when error messages are provided, since an error could be caused by different parts of the code.

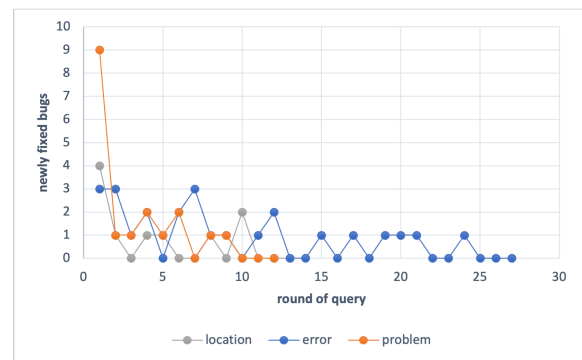


Figure 4: Number of bugs newly fixed in each round with three types of prompts

Figure 5 presents the overlapping relationship among the bugs fixed with the three types of prompts. Altogether, 34 bugs are fixed, while only 5 of them are fixed by all three prompts. Each prompt fixes 11, 7, and 2 bugs that the other two prompts fail to fix, indicating that different prompts contribute disparately to a successful bug fix.

Case Study with Problem Description. ChatGPT gains knowledge about the purpose of the program through problem descriptions, so it can find out which part of the code is inconsistent with

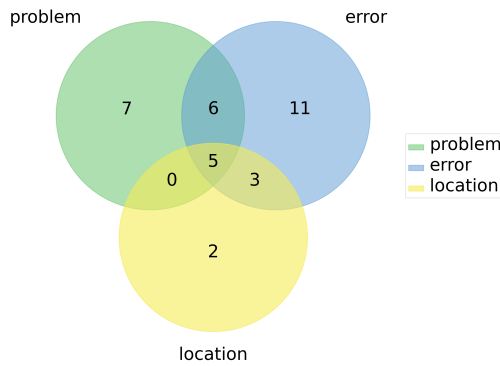


Figure 5: The overlaps of the bugs fixed by three types of prompts

```

1  ...
2  for(int i = 1; i < count; i++) {
3  int x1 = numlist.get(i - 1);
4  int x2 = numlist.get(i);
5  int dis = x2 - x1;
6  - if(dis < distance) {
7  + if(dis <= distance) {
8      System.out.println(x2);
9      break;
10 }
11 counter++;
12 }
13 ...
    
```

Listing 3: An example of fixed bug with problem description

what the coder is actually trying to do. The following is a programming problem and a buggy submission which is fixed by ChatGPT.

Problem Statement:
 Takahashi turned on a computer at time 0 and clicked the mouse N times. The i -th ($1 \leq i \leq N$) click was at time T_i .
 If he consecutively clicked the mouse at time x_1 and time x_2 (where $x_1 < x_2$), a double click is said to be fired at time x_2 if and only if $x_2 - x_1 \leq D$.
 What time was a double click fired for the first time? If no double click was fired, print -1 instead.
Input:
 The input is given from Standard Input in the following format:
 N D
 T_1 T_2 ... T_N
Output:
 If at least one double click was fired, print the time of the first such event; otherwise, print -1 .

In Listing 3, in the if statement, the operator that the problem requires is a " \leq ", as it is stated in the problem that "a double click is said to be fired at time x_2 if and only if $x_2 - x_1 \leq D$ ", but the programmer wrongly uses a " $<$ ". ChatGPT realizes the discordance between the problem and the function code, and replaces the incorrect operator to have the bug fixed.

Case Study with Error Information. ChatGPT can also infer where the bug is according to information on the type of the bug as well as on what input the program fails. Listing 4 presents a bug that can only be fixed by ChatGPT with error information. To fix the following bug, ChatGPT is informed that for the input "9 737738327422964222", the expected output is "81970925269218254" but the program actually outputs "1251275726". From the prompt,

```

1  ...
2  Scanner sc = new Scanner(System.in);
3  long A = sc.nextLong();
4  long B = sc.nextLong();
5  - int cnt = 0;
6  + long cnt = 0;
7  while(A!=B){
8      if(A>B){
9          long div = A/B;
10         A = A-B*div;
11         if(A==0){
12             div += -1;
13             cnt += div;
14             break;
15         }
16     }
17 }
    
```

Listing 4: An example of fixed bug with error information

```

1  ...
2  Scanner sc=new Scanner(System.in);
3  - long A=sc.nextInt();
4  + long A=sc.nextLong();
5  - long B=sc.nextInt();
6  + long B=sc.nextLong();
7  long sum=0;
8  while(A!=0 &&B!=0) {
9      if (A<B) {
10         long tmp=A;
11         A=B;
12         B=tmp;
13     }
14     sum+=A/B;
15     A=A%B;
16 }
    
```

Listing 5: An example of fixed bug with bug localization

ChatGPT can probably infer that there is an overflow of the output variable "cnt", so it changes the variable type from int to long.

Case Study with Bug Localization. With fault locations, ChatGPT fixes relatively fewer bugs than with the other two types of prompts. This could be attributed to the less plentiful information provided by simply buggy lines. Nevertheless, the bug position helps in some circumstances. Listing 5 presents a bug that is correctly fixed only when the bug location is offered in the prompt. While attempting to fix it based on the programming problem or the error message, ChatGPT makes the same mistake: it not only swaps the method "nextInt()" to "nextLong()" but also modifies the output "sum - 1" to "sum" at the end of the program. The prompt with the location of the bug in some ways narrows the scope of code that ChatGPT tries to edit, thus deterring ChatGPT from changing the initial correct line.

Answer to RQ2: The performance under different prompts demonstrates that, ChatGPT can benefit from more advanced prompts with additional information. For example, compared with the basic prompt, 25, 18, and 10 more bugs can be fixed with error information, problem description, and buggy lines.

4.3 RQ3: Dialogue Study

Design. During the conversation with ChatGPT, it is aware of the previous dialogues, and the response depends on both the current prompt and the context of the conversation. According to existing work [7], ChatGPT can repair more program faults through performing more dialogues. We raise this RQ to investigate the effect of dialogues on ChatGPT's repairing performance.

We focus on the unfixed bugs when giving ChatGPT either the basic prompt or the advanced prompt with error information. First, we prompt ChatGPT with the original bug and what type of error it has, and get a patch, which is then tested on the test cases. If it fails to pass any of the test cases again, we continue to conduct the next round of dialogue, in which we tell ChatGPT the failure triggered by the patch. In particular, with different types of bugs, the dialogue prompts are also different. Similar to the error information prompts described in Section 4.2, when a TLE or RE bug occurs, the [OUTPUT] will not be provided; when the bug is a CE, the dialogue will not involve any of the three elements (i.e. [INPUT], [EXPECT] and [OUTPUT]), while with a WA, all of them will be added to the prompt, shown as follows.

► Dialogue Prompt ◀

There's still a Compilation Error/ Time Limit Exceeded Error/ Runtime Error/ Wrong Answer Error in your code triggered by the input:
 ↪ [INPUT]
 The expected output is:
 ↪ [EXPECT]
 The actual output is:
 ↪ [OUTPUT]
 Try to fix it again and return the complete fix for the code.

The dialogue is continued until ChatGPT successfully repairs the bug or the round of dialogues reaches five. The API of ChatGPT has a limited length of input, so sometimes the dialogue may become too long to process. In this case, we simply delete the second round of dialogue since the first round contains the basic and essential bug information. The above process is repeated until no more new bugs are fixed for successive three times.

Results. Through dialogues, ChatGPT fixes 9 bugs among the 17 bugs that keep unfixed when only a single prompt is provided, proving that dialogues can help ChatGPT positively in repairing bugs. This can be because that although ChatGPT mistakenly fixes the bug at the first time, the knowledge towards the right fix is accumulated when we respond to a patch with what problems it has. In this way, ChatGPT can learn from its previous errors and be directed to the true location of the bug, thus the probability of generating a correct patch is improved.

There are mainly three scenarios in which ChatGPT can rectify the wrong patch with dialogues: (1) ChatGPT wrongly identifies the position of the bug, so the patch faces the same problem as what the buggy code has. After notifying ChatGPT of the problem that still exists, a new line can be located, which might be where the bug is really in. (2) ChatGPT fixes the initial problem, but there is another bug in the code that ChatGPT fails to notice. This usually happens when several bugs lie in different locations of the code. For example, in the first round of dialogue, we give a bug with a Compilation Error, caused by a syntax problem. ChatGPT corrects the syntax but a Wrong Answer problem arises. In this case, ChatGPT is told about the wrong answer and manages to find the other bug. (3) ChatGPT finds multiple suspicious code snippets, only part of which are bugs. While fixing the bugs, ChatGPT also turns some originally correct

code into the wrong one. From further dialogues, ChatGPT can realize its previous faults and produce the correct patch.

Answer to RQ3: The performance under a dialogue study demonstrates that, ChatGPT can repair more difficult-to-fix bugs with dynamic execution feedback in an interaction manner, e.g., 9 bugs that have not been fixed in previous prompts are fixed successfully.

5 DISCUSSION

In the above work, we have demonstrated that ChatGPT has an outstanding performance in program repair. However, we only ask ChatGPT to repair bugs written by human programmers in a coding contest. In this section, we aim to investigate on whether ChatGPT can fix bugs in the code generated by itself. If ChatGPT is able to fix bugs in its own code through dialogues, it will take less effort to debug manually, and help to improve humans' coding efficiency while working with ChatGPT.

To research on this question, we use the method similar to the one in RQ3, but this time the programs to be fixed by ChatGPT are given by itself. We first query ChatGPT with the problem descriptions in two AtCoder contests (i.e., abc297 and abc298) with the following prompt:

► Code Generation Prompt ◀

✎ Use Java to solve the following problem. The class name must be 'Main'. return the code in the form of markdown code block.
 ↪ [PROBLEM]

ChatGPT is asked with this prompt for 10 times and generates 10 pieces of code for each coding problem. We run them on all the test cases, and pick out the one that passes the most cases. If it passes all the test cases, it will not go into the next step. Otherwise, it is considered as the buggy code to be fixed. Next, we ask ChatGPT to fix the bug it has previously generated using the same way as in the dialogue study mentioned in Section 4.3. The dialogue round is set as 30 as we find that ChatGPT is not able to fix any bugs within a small number of dialogues, such as 5 or 10.

Among the 16 problems in the two contests, ChatGPT generates correct solutions for 3 of them when first asked to solve the problems, and then it is required to fix the remaining 13 bugs. Despite of ChatGPT's impressive bug-fixing performance in the previous experiments, it is unexpected that only 2 bugs are fixed even though we have performed 30 rounds of dialogues. This indicates that ChatGPT may have limited ability in self-repair. After careful analysis, the possible reason lies in that the edit actions to fix the developer-submitted programs are minimal (e.g., less than 6 token differences in Section 3.2), while the initial programs generated by ChatGPT are far from the correct programs and require more edits.

6 CHALLENGE AND OPPORTUNITY

Better Prompt Engineering. In the experiment, we design different prompts to feed ChatGPT with the buggy code and detailed debugging information. The results show that ChatGPT is able to fix an impressive number of bugs with the basic prompt in RQ1.

We also find advanced prompts in RQ2 and RQ3 can further boost repair performance. The quality of a prompt depends on how much information it contains and how well-crafted it is. In the bug-fixing scenario, a prompt well-crafted should contain the following element: the instruction to describe the specific task the model needs to perform, context to describe additional information that steers LLMs to better responses, and examples to describe the concrete type of the input and output samples. In the future, it is important to explore better prompts that effectively guide such LLMs to generate accurate and helpful patches for the buggy code.

Domain-specific LLMs. The investigated ChatGPT has shown outstanding performance in repairing software bugs, outperforming code domain-specific LLMs, *e.g.*, CodeT5. We think the benefits of ChatGPT lies in a wide range of training data from the internet, including books, articles, websites, and other textual data. Such diverse training data gives ChatGPT a broad understanding of language and general knowledge, while it also includes a substantial amount of information unrelated to source code and requires a huge parameter size. On the other hand, domain-specific LLMs are typically trained specifically on datasets relevant to code within a particular domain. Such domain-specific models are more focused and tailored to the code-related task and can capture domain-specific code transformation patterns, potentially resulting in more accurate patches. Overall, general-purpose LLMs (*e.g.*, ChatGPT) benefit from their broad exposure to diverse language data, enabling them to offer general programming assistance, while domain-specific LLMs (*e.g.*, CodeT5), trained solely on code-related datasets, can provide more specialized and domain-specific guidance. In the future, it is interesting to explore the advantages and disadvantages of these two types of LLMs. Besides, future work can advance the repair capabilities of LLMs by combining general language understanding with specialized domain knowledge.

ChatGPT IDE Integration. In our work, we evaluate the performance of ChatGPT with a well-constructed benchmark in terms of the number of fixed software bugs. Although ChatGPT is able to fix a considerable number of bugs, it is unclear how such LLMs perform in assisting developers in real-world development environments. Different from the existing APR pipeline that directly provides patched code, integrating LLMs like ChatGPT into IDEs can offer benefits. For example, based on the natural language and programming language understanding capabilities of ChatGPT, developers can describe the behaviors of a bug in natural language, and ChatGPT is able to analyze the buggy program and provide feasible solutions under an iterative interactive process. Besides, ChatGPT can accept more information about the bug, and even additional details or steps to reproduce the bug in IDE, enabling a more accurate diagnosis. In contrast, with limited code understanding capabilities, most existing APR techniques focus on accepting buggy code as input, leading to numerous ineffective candidate patches (*e.g.*, 1000 candidate patches per bug in CIRCLE [50]). In the future, more works are recommended to explore how LLMs streamline the typical bug-fixing workflow and offer valuable insights.

7 THREATS TO VALIDITY

The first threat to validity lies in the repair benchmarks. We construct EVALGPTFIX from AtCoder, a programming contest platform.

The collected programs are small-size algorithms, which may not accurately reflect real-world professional software repair capabilities. There is an increasing trend in using competitive programming as benchmarks in APR, as well as other tasks, such as the popular CodeNET benchmark. Besides, EVALGPTFIX contains programs with varying difficulty and types, written by programmers from diverse backgrounds. Therefore, EVALGPTFIX can effectively evaluate the repair capabilities of LLMs and foster future work on APR.

The second threat to validity comes from the compared approaches. In RQ1, we select CodeT5 and PLBART as the baselines to evaluate the effectiveness of ChatGPT. We do not consider (1) other LLMs (*e.g.*, CodeBERT [12] and GraphCodeBERT [18]) because the selected two LLMs represent state-of-the-art in bug-fixing [33]; and (2) existing APR techniques (*e.g.*, CIRCLE [50] and CoCoNut [34]) because our work focuses on LLMs on SE. However, considering that our work mainly focuses on empirical evaluations, the improvement of ChatGPT over the baselines is enough to demonstrate the promising future of boosting program repair on top of ChatGPT.

The third threat to validity is the selection of ChatGPT and APR. In our work, we regard data leakage as a common issue that may appear in a variety of SE tasks involving black-box LLMs. We only conduct experiments to evaluate the capabilities of ChatGPT in repairing software bugs. Thus, our findings may not be generalizable to other LLMs and tasks. Considering the fact that (1) ChatGPT is one of the state-of-the-art LLMs and has been extensively studied in recent works, and (2) APR plays a vital role in software development, and a number of APR works leverage LLMs to generate patches, We believe that ChatGPT and APR can indeed serve as representative examples of LLMs and SE tasks, respectively.

8 RELATED WORK

8.1 Automated Program Repair

Existing APR techniques are generally divided into traditional and learning-based ones [3, 16, 20]. Traditional APR, especially template-based APR, is proven to perform well in fixing bugs. For example, PAR [25] first proposes a patch generation method based on fix patterns, which are drawn from over 60,000 patches written manually. TBar [32] systematically summarizes frequently-used fix templates from the literature, and applies them in fixing program bugs. There are different ways to obtain fix templates. For example, In AVATAR [31], fix patterns come from code changes that can address violations in static bug detection tools. FixMiner [26] develops an automated template mining tool, in which a clustering strategy is applied to mine code changes from bugs and patches.

Recently, learning-based APR has been proposed to transform buggy code into the correct one automatically. For example, DLFix [28] raises a two-layer deep learning model to learn code transformation from bug fixes and surrounding contexts. CURE [22] pre-trains a programming language model based on a large codebase, and optimizes the searching efficiency by proposing a novel search strategy as well as using a subword tokenization technique. In this paper, we do not include these techniques as baselines because (1) previous studies [21, 53] demonstrate existing LLMs can outperform APR approaches; and (2) as an empirical study, our work mainly focuses on LLMs and select two state-of-the-art LLMs.

8.2 Large Language Models in SE

There is growing interest in leveraging LLMs in SE tasks [42, 47]. For example, Zeng et al. [58] perform an extensive study that evaluates eight pre-trained LLMs (e.g., CodeBERT [12], CodeT5 [49], and GraphCodeBERT [18]) on seven program understanding and generation tasks, and compare the pre-trained models with non-pre-trained domain-specific techniques, demonstrating that pre-trained models significantly perform better in code understanding than other state-of-the-art techniques. In the field of program repair, LLMs are exploited to directly generate patches so as to break the limitations in learning-based techniques [50, 52, 53]. AlphaRepair [53] introduces the mask prediction task of the pre-trained model CodeBERT [12] into APR, and implements a cloze-style APR tool without using any bug-fixing historic code. Prenner et al. [44] research on the program repair ability of Codex, a GPT-3-based language model aiming to translate natural language to programming language. Although Codex is not specifically trained for APR tasks, it is still effective at fixing bugs, especially those in Python language on the benchmark QuixBugs [29].

Recently, ChatGPT is attracting a lot of attention because of its stunning ability of understanding and responding to conversations started by humans. In APR, ChatGPT has been proven to have outstanding performance in fixing bugs from popular datasets (e.g. Defects4J [23] and QuixBugs [29]). Sobania et al. [45] analyze ChatGPT's program repair behavior through both giving a single request and conducting more discussions with ChatGPT. Cao et al. [7] focus on studying ChatGPT's capability in fixing deep learning programs. Xia et al. [54] propose an APR approach based on ChatGPT that fully utilizes conversations by offering instant feedbacks about previous patches. In this work, we review the data leakage issue of black-box LLMs by taking ChatGPT and APR as examples.

9 CONCLUSION

In this paper, we seek to review the overlooked data leakage issue of black-box LLMs in the SE domain. In particular, we evaluate ChatGPT's capability of program repair on a clean dataset EVALGPTFix. The results demonstrate that ChatGPT generates 109 correct patches over 151 bugs when only given the basic prompt. Besides, ChatGPT continues to fix 18, 25, and 10 additional bugs with prompts containing programming problem descriptions, error messages, and bug localizations. Through engaging in dialogues, nine more bugs are fixed by ChatGPT. These results indicate that ChatGPT has a promising bug-fixing ability, which can be further enhanced by proper prompts and more dialogues.

REFERENCES

- [1] WU Ahmad, S Chakraborty, B Ray, and KW Chang. 2021. Unified pre-training for program understanding and generation.. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- [2] Anonymous-author(s). 2023. The project website. <https://anonymous.4open.science/r/2023-EvalGPTFix/>. Lasted accessed 2023-08-01.
- [3] Johannes Bader, Andrew Scott, Michael Pradel, and Satish Chandra. 2019. Getafix: Learning to Fix Bugs Automatically. *Proceedings of the ACM on Programming Languages (OOPSLA'19)* 3, OOPSLA (2019), 1–27.
- [4] Samuel Benton, Xia Li, Yiling Lou, and Lingming Zhang. 2020. On the Effectiveness of Unified Debugging: An Extensive Study on 16 Program Repair Systems. In *2020 35th IEEE/ACM International Conference on Automated Software Engineering (ASE'20)*. IEEE, 907–918.
- [5] Tom Britton, Lisa Jeng, Graham Carver, and Paul Cheak. 2013. Reversible Debugging Software “quantify the Time and Cost Saved Using Reversible Debuggers”. (2013).
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language Models Are Few-shot Learners. In *Proceedings of the Advances in Neural Information Processing Systems*, Vol. 33. 1877–1901.
- [7] Jialun Cao, Meiziniu Li, Ming Wen, and Shing-chi Cheung. 2023. A study on Prompt Design, Advantages and Limitations of ChatGPT for Deep Learning Program Repair. *arXiv preprint arXiv:2304.08191* (2023).
- [8] Zimin Chen, Steve James Kommrusch, and Martin Monperrus. 2022. Neural Transfer Learning for Repairing Security Vulnerabilities in C Code. *IEEE Transactions on Software Engineering (TSE)* (2022).
- [9] Jianlei Chi, Yu Qu, Ting Liu, Qinghua Zheng, and Heng Yin. 2022. Seqtrans: Automatic Vulnerability Fix Via Sequence to Sequence Learning. *IEEE Transactions on Software Engineering (TSE)* (2022).
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [11] Thomas Durieux and Martin Monperrus. 2016. Dynamoth: Dynamic Code Synthesis for Automatic Program Repair. In *Proceedings of the 11th International Workshop on Automation of Software Test (AST'16)*. 85–91.
- [12] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiao Cheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. Codebert: A Pre-trained Model for Programming and Natural Languages. In *Findings of the Association for Computational Linguistics (EMNLP'20)*. 1536–1547.
- [13] Michael Fu, Chakkrit Tantithamthavorn, Trung Le, Van Nguyen, and Phung Dinh. 2022. Vulrepair: A T5-based Automated Software Vulnerability Repair. In *the ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE'22)*.
- [14] Xiang Gao, Yannic Noller, and Abhik Roychoudhury. 2022. Program Repair. *arXiv preprint arXiv:2211.12787* (2022).
- [15] Luca Gazzola, Daniela Micucci, and Leonardo Mariani. 2019. Automatic Software Repair: A Survey. *IEEE Transactions on Software Engineering (TSE)* 45, 1 (2019), 34–67.
- [16] Ali Ghanbari, Samuel Benton, and Lingming Zhang. 2019. Practical Program Repair Via Bytecode Mutation. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA'19)*. 19–30.
- [17] Daya Guo, Shuai Lu, Nan Duan, Yanlin Wang, Ming Zhou, and Jian Yin. 2022. Unixcoder: Unified Cross-modal Pre-training for Code Representation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL'22)*. 7212–7225.
- [18] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, et al. 2021. Graphcodebert: Pre-training Code Representations with Data Flow. In *Proceedings of the 9th International Conference on Learning Representations (ICLR'21)*. 1–18.
- [19] Md Mahim Anjum Haque, Wasi Uddin Ahmad, Ismini Lourentzou, and Chris Brown. 2022. Fixeval: Execution-based Evaluation of Program Fixes for Programming Problems. *arXiv preprint arXiv:2206.07796* (2022).
- [20] Jiajun Jiang, Yingfei Xiong, Hongyu Zhang, Qing Gao, and Xiangqun Chen. 2018. Shaping Program Repair Space with Existing Patches and Similar Code. In *Proceedings of the 27th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA'18)*. 298–309.
- [21] Nan Jiang, Kevin Liu, Thibaud Lutellier, and Lin Tan. 2023. Impact of code language models on automated program repair. In *Proceedings of the 45th International Conference on Software Engineering (ICSE'23)*. 1430–1442.
- [22] Nan Jiang, Thibaud Lutellier, and Lin Tan. 2021. Cure: Code-aware Neural Machine Translation for Automatic Program Repair. In *Proceedings of the 43rd IEEE/ACM International Conference on Software Engineering (ICSE'21)*. 1161–1173.
- [23] René Just, Darioush Jalali, and Michael D Ernst. 2014. Defects4j: A Database of Existing Faults to Enable Controlled Testing Studies for Java Programs. In *Proceedings of the 23rd International Symposium on Software Testing and Analysis (ISSTA'14)*. 437–440.
- [24] Anjan Karmakar, Julian Aron Prenner, Marco D'Ambros, and Romain Robbes. 2022. Codex Hacks HackerRank: Memorization Issues and a Framework for Code Synthesis Evaluation. *arXiv preprint arXiv:2212.02684* (2022).
- [25] Dongsun Kim, Jaechang Nam, Jaewoo Song, and Sunghun Kim. 2013. Automatic patch generation learned from human-written patches. In *2013 35th International Conference on Software Engineering (ICSE'13)*. IEEE, 802–811.
- [26] Anil Koyuncu, Kui Liu, Tegawendé F Bissyandé, Dongsun Kim, Jacques Klein, Martin Monperrus, and Yves Le Traon. 2020. Fixminer: Mining Relevant Fix Patterns for Automated Program Repair. *Empirical Software Engineering (ESE)* 25, 3 (2020), 1980–2024.
- [27] Claire Le Goues, ThanhVu Nguyen, Stephanie Forrest, and Westley Weimer. 2012. GenPro: A Generic Method for Automatic Software Repair. *IEEE Transactions on Software Engineering (TSE)* 38, 01 (2012), 54–72.
- [28] Yi Li, Shaohua Wang, and Tien N Nguyen. 2020. Dlfix: Context-based Code Transformation Learning for Automated Program Repair. In *Proceedings of the*

- 42nd ACM/IEEE International Conference on Software Engineering (ICSE'20). 602–614.
- [29] Derrick Lin, James Koppel, Angela Chen, and Armando Solar-Lezama. 2017. Quixbugs: A Multi-lingual Program Repair Benchmark Set Based on the Quixey Challenge. In *Proceedings Companion of the 2017 ACM SIGPLAN International Conference on Systems, Programming, Languages, and Applications: Software for Humanity (SPLASH Companion '17)*. 55–56.
- [30] Chao Liu, Xuanlin Bao, Hongyu Zhang, Neng Zhang, Haibo Hu, Xiaohong Zhang, and Meng Yan. 2023. Improving ChatGPT Prompt for Code Generation. *arXiv preprint arXiv:2305.08360* (2023).
- [31] Kui Liu, Anil Koyuncu, Dongsun Kim, and Tegawendé F Bissyandé. 2019. Avatar: Fixing Semantic Bugs with Fix Patterns of Static Analysis Violations. In *Proceedings of the 26th IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER'19)*. 1–12.
- [32] Kui Liu, Anil Koyuncu, Dongsun Kim, and Tegawendé F Bissyandé. 2019. Tbar: Revisiting Template-based Automated Program Repair. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA'19)*. 31–42.
- [33] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. 2021. Codexglue: A Machine Learning Benchmark Dataset for Code Understanding and Generation. *arXiv preprint arXiv:2102.04664* (2021).
- [34] Thibaud Lutellier, Hung Viet Pham, Lawrence Pang, Yitong Li, Moshi Wei, and Lin Tan. 2020. Coconut: Combining Context-aware Neural Translation Models Using Ensemble for Program Repair. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA'20)*. 101–114.
- [35] Matias Martinez and Martin Monperrus. 2016. Astor: A Program Repair Library for Java. In *Proceedings of the 25th International Symposium on Software Testing and Analysis (ISSTA'16)*. 441–444.
- [36] Matias Martinez and Martin Monperrus. 2018. Ultra-large Repair Search Space with Automatically Mined Templates: The Cardumen Mode of Astor. In *Proceedings of the International Symposium on Search Based Software Engineering (SSBSE'18)*. Springer, 65–86.
- [37] Antonio Mastropaolo, Simone Scalabrino, Nathan Cooper, David Nader Palacio, Denys Poshyvanyk, Rocco Oliveto, and Gabriele Bavota. 2021. Studying the Usage of Text-to-text Transfer Transformer to Support Code-related Tasks. In *Proceedings of the 43rd IEEE/ACM International Conference on Software Engineering (ICSE'21)*. 336–347.
- [38] Sergey Mechtaev, Jooyong Yi, and Abhik Roychoudhury. 2016. Angelix: Scalable Multiline Program Patch Synthesis Via Symbolic Analysis. In *Proceedings of the 38th international conference on software engineering (ICSE'16)*. 691–701.
- [39] Martin Monperrus. 2018. Automatic Software Repair: A Bibliography. *ACM Computing Surveys (CSUR)* 51, 1 (2018), 1–24.
- [40] Martin Monperrus. 2022. The Living Review on Automated Program Repair. (2022).
- [41] Noor Nashid, Mifta Sintaha, and Ali Mesbah. 2023. Retrieval-based prompt selection for code-related few-shot learning. In *Proceedings of the 45th International Conference on Software Engineering (ICSE'23)*.
- [42] Changan Niu, Chuanyi Li, Vincent Ng, Jidong Ge, Liguang Huang, and Bin Luo. 2022. Spt-code: Sequence-to-sequence Pre-training for Learning the Representation of Source Code. In *Proceedings of the 44th International Conference on Software Engineering (ICSE'22)*. 2006–2018.
- [43] OpenAI. 2023. ChatGPT: Optimizing Language Models for Dialogue. <https://openai.com/blog/chatgpt>. Lasted accessed 2023-08-01.
- [44] Julian Aron Prenner, Hlib Babii, and Romain Robbes. 2022. Can Openai's Codex Fix Bugs? An Evaluation on Quixbugs. In *Proceedings of the Third International Workshop on Automated Program Repair (APR'22)*. 69–75.
- [45] Dominik Sobania, Martin Briesch, Carol Hanna, and Justyna Petke. 2023. An analysis of the automatic bug fixing performance of chatgpt. In *Proceedings of the Fourth International Workshop on Automated Program Repair (APR'23)*. 1–8.
- [46] Weisong Sun, Chunrong Fang, Yudu You, Yun Miao, Yi Liu, Yuekang Li, Gelei Deng, Shenghan Huang, Yuchen Chen, Quanjun Zhang, et al. 2023. Automatic Code Summarization via ChatGPT: How Far Are We? *arXiv preprint arXiv:2305.12865* (2023).
- [47] Michele Tufano, Jevgenija Pantuchina, Cody Watson, Gabriele Bavota, and Denys Poshyvanyk. 2019. On Learning Meaningful Code Changes Via Neural Machine Translation. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE'19)*. IEEE, 25–36.
- [48] Michele Tufano, Cody Watson, Gabriele Bavota, Massimiliano Di Penta, Martin White, and Denys Poshyvanyk. 2019. An Empirical Study on Learning Bug-fixing Patches in the Wild Via Neural Machine Translation. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 28, 4 (2019), 1–29.
- [49] Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. 2021. Codet5: Identifier-aware Unified Pre-trained Encoder-decoder Models for Code Understanding and Generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP'21)*. 8696–8708.
- [50] Yuan Wei, Zhang Quanjun, He Tiek, Fang Chunrong, Hung Nguyen Quoc Viet, Hao Xiaodong, and Yin Hongzhi. 2022. Circle: Continual Repair across Programming Languages. In *Proceedings of the 31th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA'22)*. 427–438.
- [51] W. E. Wong, R. Gao, Y. Li, R. Abreu, and F. Wotawa. 2016. A Survey on Software Fault Localization. *IEEE Transactions on Software Engineering (TSE)* 42, 8 (Aug. 2016), 707–740.
- [52] Chunqiu Steven Xia, Yuxiang Wei, and Lingming Zhang. 2023. Automated Program Repair in the Era of Large Pre-trained Language Models. In *Proceedings of the 45th International Conference on Software Engineering (ICSE'23)*. 1482–1494.
- [53] Chunqiu Steven Xia and Lingming Zhang. 2022. Less Training, More Repairing Please: Revisiting Automated Program Repair Via Zero-shot Learning. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE'22)*. 959–971.
- [54] Chunqiu Steven Xia and Lingming Zhang. 2023. Keep the Conversation Going: Fixing 162 out of 337 bugs for \$0.42 each using ChatGPT. *arXiv preprint arXiv:2304.00385* (2023).
- [55] Yingfei Xiong, Jie Wang, Runfa Yan, Jiachen Zhang, Shi Han, Gang Huang, and Lu Zhang. 2017. Precise Condition Synthesis for Program Repair. In *Proceedings of the 39th IEEE/ACM International Conference on Software Engineering (ICSE'17)*. IEEE, 416–426.
- [56] Yuan Yuan and Wolfgang Banzhaf. 2018. Arja: Automated Repair of Java Programs Via Multi-objective Genetic Programming. *IEEE Transactions on Software Engineering (TSE)* 46, 10 (2018), 1040–1067.
- [57] Zhiqiang Yuan, Yiling Lou, Mingwei Liu, Shiji Ding, Kaixin Wang, Yixuan Chen, and Xin Peng. 2023. No More Manual Tests? Evaluating and Improving ChatGPT for Unit Test Generation. *arXiv preprint arXiv:2305.04207* (2023).
- [58] Zhengran Zeng, Hanzhuo Tan, Haotian Zhang, Jing Li, Yuqun Zhang, and Lingming Zhang. 2022. An extensive study on pre-trained models for program understanding and generation. In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis*. 39–51.
- [59] Quanjun Zhang, Chunrong Fang, Weisong Sun, Yan Liu, Tiek He, Xiaodong Hao, and Zhenyu Chen. 2023. Boosting Automated Patch Correctness Prediction via Pre-trained Language Model. *arXiv preprint arXiv:2301.12453* (2023).
- [60] Qihao Zhu, Zeyu Sun, Wenjie Zhang, Yingfei Xiong, and Lu Zhang. 2023. Tare: Type-Aware Neural Program Repair. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 1443–1455.