Integrating Various Software Artifacts for Better LLM-based Bug Localization and Program Repair

QIONG FENG, Nanjing University of Science and Technology, China XIAOTIAN MA, Nanjing University of Science and Technology, China JIAYI SHENG, Nanjing University of Science and Technology, China ZIYUAN FENG, Nanjing University of Science and Technology, China WEI SONG, Nanjing University of Science and Technology, China PENG LIANG, Wuhan University, China

LLMs have garnered considerable attention for their potential to streamline Automated Program Repair (APR). LLM-based approaches can either insert the correct code using an infilling-style technique or directly generate patches when provided with buggy methods, aiming for plausible patches to pass all tests. However, most of LLM-based APR methods rely on a single type of software information, such as issue descriptions or error stack traces, without fully leveraging a combination of diverse software artifacts. Human developers, in contrast, often use a range of information — such as debugging data, issue discussions, and error stack traces — to diagnose and fix bugs. Despite this, many LLM-based approaches do not explore which specific types of software information best assist in localizing and repairing software bugs. Addressing this gap is crucial for advancing LLM-based APR techniques.

To investigate this and mimic the way human developers fix bugs, we propose DEVLORE (short for DEV eloper Localization and Repair). In this framework, LLMs first use issue content (description and discussion) and stack error traces to localize buggy methods, then rely on debug information in buggy methods and issue content and stack error to localize buggy lines and generate valid patches. We evaluated the effectiveness of issue content, error stack traces, and debugging information in bug localization and automatic program repair. Our results show that while issue content and error stack is particularly effective in assisting LLMs with fault localization and program repair respectively, different types of software artifacts complement each other in addressing various bugs. By incorporating these three types of artifacts and using the Defects4J v2.0 dataset for evaluation, DEVLORE successfully localizes 49.3% of single-method bugs and generates 56.0% plausible patches. Additionally, DEVLORE can localize 47.6% of non-single-method bugs and generates 14.5% plausible patches. Moreover, our framework streamlines the end-to-end process from buggy source code to a complete repair, and achieves a 39.7% and 17.1% of single-method and non-single-method bug repair rate, outperforming current state-of-the-art APR methods. The source code and experimental results of this work for replication are available at https://github.com/XYZboom/DEVLoRe.

 ${\tt CCS\ Concepts: \bullet Software\ and\ its\ engineering \to Software\ maintenance\ tools; Software\ verification\ and\ validation.}$

Additional Key Words and Phrases: Large Language Model, Automatic Program Repair, Fault Localization, A

Authors' Contact Information: Qiong Feng, Nanjing University of Science and Technology, Nanjing, China, qiongfeng@njust.edu.cn; Xiaotian Ma, Nanjing University of Science and Technology, Nanjing, China, xyzboom@njust.edu.cn; Jiayi Sheng, Nanjing University of Science and Technology, Nanjing, China, shengjiayi@njust.edu.cn; Ziyuan Feng, Nanjing University of Science and Technology, Nanjing, China, azumaseren@njust.edu.cn; Wei Song, Nanjing University of Science and Technology, Nanjing, China, wsong@njust.edu.cn; Peng Liang, Wuhan University, Wuhan, China, liangp@whu.edu.cn.

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 the author(s) 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.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM 1557-7392/2024/3-ART

https://doi.org/XXXXXXXXXXXXXXX

ACM Reference Format:

1 Introduction

Automatic Program Repair (APR) streamlines the process of identifying and correcting code defects, significantly reducing the time and effort required for manual bug fixing [18]. Traditional APR techniques employ various methods, including template-based approaches [7, 15, 17, 22, 24, 26] and neural machine translation (NMT) [5, 10, 21, 46], to generate potential patches that are both syntactically valid and semantically meaningful. Although these methods can produce correct patches for certain bugs, they also have notable limitations. NMT models rely heavily on bug-fixing training data, making them unable to generate patches for new or unseen types of bugs. On the other hand, template-based approaches suffer from a limited set of templates and struggle to address more complex, nontrivial bug fixes [38].

Recent studies have explored the use of LLMs for APR, either by having LLMs fill in the correct code in buggy methods using an infilling-style technique or by directly generating patches when provided with buggy methods [9, 20, 31, 34–39, 44]. The initial results demonstrate the ability of LLMs to correctly repair real-world bugs, including those that were previously unrepairable by existing APR approaches. One of the best-performing frameworks is Agentless [35], which feeds an LLM with only issue description and can automatically solve GitHub issues. These promising outcomes highlight the potential of LLMs to develop more effective APR methods.

However, there are still two major **limitations** that need to be addressed:

- Current LLM-based approaches do not fully incorporate various types of software artifacts.
 Most LLMs rely on just one or two kinds of software artifacts, such as issue descriptions or
 code structure, while other important artifacts are underutilized. For instance, debugging
 information, which is a critical tool for human developers in diagnosing and resolving bugs,
 is often not fully leveraged.
- While various LLM-based approaches make use of different software artifacts, such as issue
 descriptions and error stack traces, it remains unclear which specific type of information
 most effectively aids LLMs in localizing and automatically repairing software bugs.

To address these two limitations and further explore the ability of LLMs to localize and fix software bugs (bug, defect, and fault are used interchangeably in this paper), we propose feeding LLMs different types of software artifacts to determine which information best leverages their capabilities in bug localizing and fixing. The rationale behind this approach is that human developers typically do not rely on a single type of information when localizing and fixing software bugs. Instead, they combine various sources of information in software development, such as issue descriptions, proof of concept (PoC), stack traces, discussions in issues, and more. Based on developers' experience, having more information helps to better understand the root cause of bugs, ultimately leading to more effective bug localization and fixes.

To achieve this, we propose the DEVLORE framework, which asks LLMs to mimic human developers for bug localization and program repair. Along with this framework, we design two tools to extract executed methods in failed test cases and debugging information of buggy methods. In this framework, we feed the chosen LLM with three types of software artifacts: issue content (including issue description and discussion), error stack trace, and debug information. Then, using the well-known Defects4J dataset [13], we evaluate how these different types of software artifacts contribute to effectively localizing and fixing bugs. Our experiment results demonstrate that issue

content is the most effective indicator for buggy method's localization, achieving 43.6% for localizing single-method bugs and 40.6% for localizing non-single-method bugs. Stack error trace proves to be the best indicator for single-method bug fixing, with a precision of 27.0%. Additionally, our findings highlight that different types of software information can complement each other in the process of localizing and fixing bugs. By combining issue content and error stack trace, we achieve a state-of-the-art performance of 49.3% in localizing the single-method bugs. Furthermore, incorporating issue content, error stack trace, and debugging information results in a fix rate of 43.1% for single-method bugs.

As our DEVLoRE framework does not specify that the localization and fix of buggy methods should be singular, it can localize and fix bugs across multiple methods (non-single-method bugs). For example, by combining all three software artifacts, DEVLoRE can fix 9.4% of non-single-method bugs with provided buggy locations. Moreover, similar to Agentless [35], DEVLoRE relies on LLMs throughout the entire end-to-end process of fault localization and program repair. It can successfully fix 28.0% of single-method bugs and 11.2% of non-single-method bugs when combining all three artifacts. To the best of our knowledge, our approach outperforms current state-of-the-art methods.

To summarize, our **contributions** in this paper are as follows:

- (1) To the best of our knowledge, this work is the first to compare different software artifacts in assisting LLMs' ability to perform fault localization and program repair. The results can help developers understand LLMs' potential when provided with a variety of software artifacts.
- (2) We propose a simple and lightweight framework that leverages LLMs to conduct an end-to-end process for bug localization and program repair. Accompanying this framework we design a strict input/output prompt. The DEVLORE framework demonstrates a strong ability to localize and fix more software bugs in less time and at a lower cost, compared to current state-of-the-art methods.

The paper is structured as follows: Section 2 discusses the motivation of using various software artifacts, Section 3 outlines the proposed approach, Section 4 introduces the experiment setup and the research questions, Sections 5 and 6 present and discuss the results respectively, Section 7 reviews the related work, and Section 8 concludes this work with future directions.

2 Three Motivating Examples

Our assumption is that LLM-based program repair tools have the potential to be much more effective if they can leverage a variety of software artifacts, such as code snippets, version history, documentation, and even testing outputs. By providing LLMs with a rich set of contextual data, they can understand and localize bugs more accurately and offer more precise fixes. This assumption is based on a human software engineer's daily practice. When a human developer tries to fix a bug, they would examine various resources such as the issue descriptions, error stack, debugging information, and test cases, until a solution is identified. However, currently most LLM-based repair approaches underutilize or use limited software artifacts. Here, we provide three examples to demonstrate the motivation of this work and the ability of LLMs to localize and fix bugs when given access to different types of software artifacts.

2.1 Debugging Information to the Rescue

The Lang 1b bug involves the *createNumber* method, which takes a string as its parameter and returns a number represented by that string. While studying the Lang 1b bug, we noticed that the newly added failing test case is the input "0x80000000" for the *createNumber* method. When this test case is provided, the buggy *createNumber* method incorrectly treats the input string as an

integer, when it should be a long integer instead. The code snippet in Figure 1 shows a portion of the buggy *createNumber* method in Lang 1b. The local variable hexDigits represents the valid hexadecimal digits. As human software developers, we can quickly recognize that at Line 471, the condition should be " $hexDigits > 8 \parallel (hexDigits == 8 \&\& firstSigDigit > '7')$ " (firstSigDigit refers to the first valid hexadecimal digit).

```
if (pfxLen > 0) { // we have a hex number
466
                     final int hexDigits = str.length() - pfxLen;
467
468
                     if (hexDigits > 16) { // too many for Long
                         return createBigInteger(str);
469
470
471
                     if (hexDigits > 8) { // too many for an int
472
                         return createLong(str);
473
                     }
474
                     return createInteger(str);
475
```

Fig. 1. Code snippet for Lang 1b

For this bug, we provided the LLM with Lang 1b's issue content (https://issues.apache.org/jira/browse/LANG-747) and the stack error message shown in Figure 2, but both could not help the LLM generate a plausible fix or near-to-correct patch. This is because the leap required in reasoning is significant, and the buggy method is quite long. Neither of the information provided above could pinpoint the buggy line precisely, let alone generate a plausible fix. To address this, we introduced debugging information to help the LLM understand the necessary changes. One piece of debug information we provided was a series of variable-value pairs, commons.lang3.math.NumberUtils:createNumber:468 [hexDigits:8]. This indicates that the code is about to execute Line 468 and that the value of hexDigits

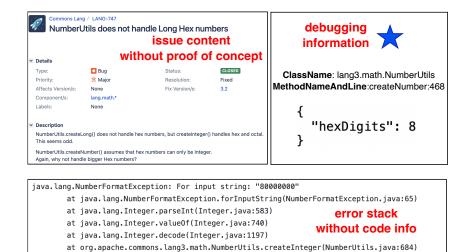


Fig. 2. Issue content, error stack and debug info of Lang 1b

at org.apache.commons.lang3.math.NumberUtils.createNumber(NumberUtils.java:474) at org.apache.commons.lang3.math.NumberUtilsTest.TestLang747(NumberUtilsTest.java:256)

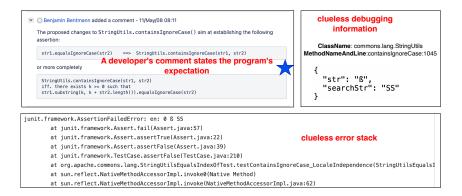


Fig. 3. Issue content, error stack and debug info for Lang 40b

is 8 for the given test case. This enables the LLM to understand why the code on Line 474 was executed instead of Line 471. In one of the responses, the LLM suggested changing Line 468 to *if* ($hexDigits > 16 \parallel (hexDigits == 8 \&\& str.charAt(2) >= '8')$), which, although not completely accurate, demonstrates that the provided debug information has a positive impact on fixing the bug.

2.2 Issue Content to the Rescue

Sometimes, the discussions and Proof of Concept (PoC) included in the issue content can help the LLM better understand the expectations for bug repair. Figure 3 shows a developer's comment in the issue content, which states the expected behavior of the program (https://issues.apache.org/jira/browse/LANG-432). The original code attempted to convert both input strings to uppercase using String.toUpperCase() and then return the result of invoking the *contains* method. However, String.toUpperCase() is locale-sensitive, which makes it unsuitable for case-insensitive comparisons. For example, the character **0x00DF** represents "ß" in Unicode. If we apply String.toUpperCase() to this character, it becomes "SS" in the Turkey locale. Consequently, comparing "ss" with "ß" would result in an equal match. Therefore, in this case, we should not use String.toUpperCase() but instead compare the characters individually.

The debug information, org.apache.commons.lang.StringUtils:containsIgnoreCase:1045 {"str":" β ", "searchStr":"SS"}, and the error stack trace shown in Figure 3 merely re-display the test cases and do not clearly highlight the relationship between the character 0x00DF and 'SS'. As a result, neither the error stack trace nor the debug information provides any useful clues about the bug, and LLMs cannot generate a plausible fix when provided with either the stack message or the debug information alone. However, with the description and expected behavior outlined in the issue content, GPT-40 is able to generate a correct patch.

2.3 Error Stack to the Rescue

When a bug triggers an exception, as opposed to merely being an unexpected behavior, the error stack trace generated at the point where the exception occurs can be helpful in diagnosing and fixing the error. Stack trace provides a detailed record of the sequence of method calls that led to an exception, making it easier to pinpoint the exact location in the code where the bug arises.

For example, in Lang 39b shown in Figure 4, the developer forgot to check for *null* values in the elements of the input parameters *replacementList* and *searchList*, which can lead to a NullPointerException (NPE). However, if we do not inform the LLM about the occurrence of an NPE in the error stack, the LLM will not be able to identify where *null* detection is necessary. The stack trace

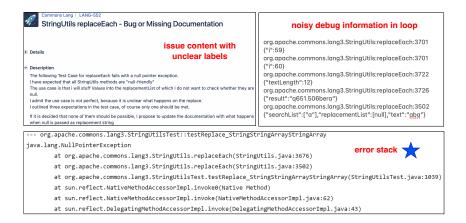


Fig. 4. Issue content, error stack and debug info for Lang 39b

```
diff --git a/src/java/org/apache/commons/lang3/StringUtils.java b/src/java/org/apache/commons/lang3/StringUtils.java
index 14563aa6..2d12a357 100644
--- a/src/java/org/apache/commons/lang3/StringUtils.java
+++ b/src/java/org/apache/commons/lang3/StringUtils.java
@@ -3673,7 +3673,8 @@ public class StringUtils {

    // count the replacement text elements that are larger than their corresponding text being replaced
    for (int i = 0; i < searchList.length; i++) {

        int greater = replacementList[i].length() - searchList[i].length();

        int greater = (replacementList[i] != null ? replacementList[i].length() : 0) - searchList[i].length();

        if (greater > 0) {
              increase += 3 * greater; // assume 3 matches
        }
}
```

Fig. 5. A patch generated by GPT-4o-mini for Lang 39b

clearly indicates that the NPE occurs on Line 3676 in the class *StringUtils*. By providing this stack information to the LLM, it is able to generate a repair, as shown in Figure 5.

2.4 Observations From the Above Three Examples

From these three examples, we can see that issue content, error stack traces, and debugging information can complement each other in bug detection and fixing. Each type of artifact provides a different level of context, and when combined, they create a more comprehensive understanding of the bug. Issue content offers the high-level description and expectations for the fix, error stack traces pinpoint the exact location of the failure, and debugging information provides granular details about the variable states and flow of execution. By incorporating these diverse artifacts, the LLM's repair process becomes more holistic.

3 Approach

As shown in Figure 6, our approach utilizes LLMs in two distinct steps: bug localization and program repair. Along with these two steps, we developed two dynamic tools — MethodRecorder and DebugRecorder — to assist LLMs in handling bug localization and fix. MethodRecorder tracks the methods executed during the failing test case(s), while DebugRecorder is designed to extract debug information from the buggy method(s). Debug information is stored in a list, where each

Table 1. Prompts used in localizing buggy methods, buggy lines and generating patches

GENERAL TASK PROMPT	Input Prompt	EXPECTED OUTPUT PROMPT
You are a Software Engineer. Review the following skeleton of classes, test case(s), and exception that occurs when doing the test. Provide a set of locations that need to be edited to fix the bug. The locations must be specified as method names or field names.	### Skeleton of Classes ### {related_methods} {error_stack} {issue_content}	Please localize class name and method names or field names that need to be edited. Examples: path.to.ClassA::methodA path.to.ClassA::methodB path.to.ClassB::methodA
You are a Software Engineer. Review the following skeleton of classes, test case(s), and exception that occurs when doing the test. Provide a set of locations that need to be edited to fix the bug. The locations must be specified as line number in class.	### Skeleton of Classes ### {buggy_methods} {error_stack} {issue_content} {debugging_info}	Please localize class name and line number that need to be edited. Examples: path.to.ClassA line: 20 line: 45 line: 46 line: 47
You are a Software Engineer. Review the following methods and(or) fields of classes, test case(s), and exception that occurs when doing the test. Try to fix the bug.	### Skeleton of Classes ### {buggy_methods} ### Possible bug locations (for your reference only) ### {buggy_lines} {error_stack}	Please generate *SEARCH/REPLACE* edits to fix the bug based on the info given above. Every *SEARCH/REPLACE* edit must use this format: 1. The file path 2. The start of search block: <<<<< SEARCH 3. A contiguous chunk of lines to search in the existing source code 4. The dividing line: ======= 5. The lines to replace into the source code

entry includes the method name, line number, variable names, and the JSON-formatted content of these variables. Table 1 shows the prompts used in our process. The first, second, and third rows correspond to the prompts used for the tasks of localizing buggy methods, localizing buggy lines, and generating patches, respectively. The first two rows of prompts are used for the bug localization task and the last row of prompt is used for the program repair task. For each task, we provide the LLM with <General Task Prompt> + <Input Prompt> + <Expected Output Prompt> to receive the response from the LLM.

PLACE

{issue content}

{debugging info}

to fix the bug.

5. The lines to replace into the source code

6. The end of the replace block: >>>>> RE-

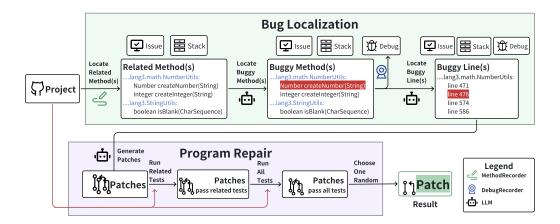


Fig. 6. Our approach for bug localization and program repair

3.1 Localize Buggy Method(s)

We emulate the behavior of human developers, who use errors in stack traces and proof of concept (PoC) from issue content to localize buggy methods. To this end, we apply MethodRecorder to trace the methods invoked during the execution of failing test cases. We then collect the signatures of these invoked methods and prompt the LLM to identify the buggy method(s). In this step, we also provide the LLM with any available error stack traces and issue content, including general description of the issue and developers' discussion/comment under this issue. We choose not to use debugging information because it is typically gathered at runtime and involves instrumentation (see <code>java.lang.instrument.Instrumentation</code> [1]), which can be computationally expensive. In developers' daily practice, they set just a few of breakpoints and collect important information about variables. But setting breakpoints automatically is not practical. Therefore, before localizing buggy method(s), minimizing instrumentation and debugging information is more efficient.

3.2 Localize Buggy Line(s)

Once the buggy method(s) have been localized in the previous step, we utilize DebugRecorder to capture detailed information about variable names and their values within the identified buggy method(s). The rationale for introducing debugging information at this stage is that it provides a more granular view, which is particularly beneficial for pinpointing the specific lines responsible for the bug and for understanding the internal state of variables. The debugging information with variables' value is often unnecessary during the initial localization step, but becomes crucial when the focus shifts to identifying precise faults within a known problematic method.

In this process, we gather debugging information as a list of variable-value pairs that reflect the state at specific lines of code within the buggy method(s). This allows the LLM to understand the conditions and data flows that may contribute to the bug. Combined with the issue content, such as descriptions and discussions, along with stack traces, this debugging information is fed into the LLM along with the body of the buggy method(s). Together, these resources assist the LLM in localizing the exact line(s) within the method(s) where the bug manifests, providing the context needed to understand and resolve the bug effectively.

3.3 Patch Generation

For patch generation, we provide all relevant information to maximize the LLM's ability to produce an accurate repair. This includes the localized buggy buggy line(s) and method(s) identified in previous steps, as well as the complete body of the buggy method(s). Additionally, we provide supplementary materials such as the issue content, which may contain descriptions and discussions about the bug, the stack trace from the error, and detailed debugging information captured at runtime. By supplying this comprehensive context, we enable the LLM to better understand the specific code behavior that led to the error and to use this knowledge to generate a more targeted and effective patch for the bug. For each bug, we use the LLM to generate multiple potential patches in *.diff* format, which can be directly applied to the buggy code.

3.4 Patch Validation

We apply the patches generated in the previous step to the buggy code and test them to evaluate their effectiveness. If the patch passes the initial failed tests, we proceed it to pass all tests in case that it incur any regression. If it passed all the tests, then it is considered a plausible fix. To facilitate this process, we utilize the Defects4J framework to run the tests and verify the patches.

To be consistent with existing studies [20, 47], if the patch is semantically equivalent to the original patch provided by developers, it is considered a correct patch. As part of this validation process, the patch undergoes manual inspection and cross-checking by two experienced developers to ensure its correctness. This step is crucial to verify that the patch not only passes all automated tests but also aligns with the intended behavior from a developer's perspective. The manual inspection by developers serves as a final quality control step to ensure that the patch addresses the bug correctly and does not introduce new bugs.

4 Experiment Setup

4.1 Dataset

We used the well-established Defects4J benchmark [13] for our experiments, specifically leveraging both version 1.2 and version 2.0. Defects4J v1.2 contains 391 bugs from six real-world projects (note that there are 395 bugs in v1.2, but due to the update to Java 8, four bugs can no longer be reproduced), while v2.0 includes an additional 444 bugs from 11 real-world projects. We chose Defects4J benchmark for two main reasons. First, it offers a diverse set of software artifacts — such as issue URLs, error stacks, and test cases — that are well-suited for our proposed ARP framework and help maximize its capabilities. Second, Defects4J benchmark has been widely used in prior research, which facilitates a fair and meaningful comparison of our approach.

4.2 Parameters in Experiments

Due to the extremely long debugging information for some bugs recorded by DebugRecorder, along with issue content and error stack, which exceeds the token limit of most LLMs, we opted for the GPT-4 series model, known for its ability to handle long token sequences with a 128k token context window. To manage costs, we selected the cost-effective GPT-4o-mini model, also with a 128k token context window, and used the default settings (*Temperature*=0.5, Top-*p*=1). We first conducted 10 manual experiments in each step (localizing buggy methods, localizing buggy lines, and generating patches) to determine how many times the LLM should be called at each stage. The rationale is that calling the LLM more times does not significantly improve performance beyond a certain point. Based on these manual experiments, we finalized the settings of the parameters. For buggy method localization, we call the LLM only once. For buggy line localization, we call the LLM 10 times, collecting 10 responses, each identifying potential buggy lines. We then filter out

duplicate responses (i.e., identical buggy line locations). Finally, for each unique response, we ask the LLM to generate patches three times, ensuring that multiple potential fixes are considered.

4.3 Research Questions

In this work, we aim to answer the following research questions (RQs) for evaluating the effectiveness of various software artifacts in assisting the chosen LLM in bug localization and APR task.

- RQ1. Which software artifacts can better assist LLMs in localizing buggy methods and buggy lines when provided with source code? This RQ aims to identify the specific types of software artifacts - issue content, debugging information and error stack trace that enhance the ability of LLMs to accurately pinpoint methods and code that contain bugs.
- RQ2. Which software artifacts can better assist LLMs in generating plausible patches
 when provided with buggy methods? By answering this RQ, we can identify the specific
 types of software artifacts in assisting LLMs to produce plausible patches which can pass
 all unit tests.
- RQ3. What is the overall performance of DEVLoRE in bug localization and program repair? Unlike some LLM-based approaches that focus on either fault localization or program repair while relying on traditional methods such as spectrum-based approaches [3] for fault localization [20, 36], DEVLoRE supports an end-to-end bug identification and program repair workflow. This RQ aims to evaluate the overall performance of DEVLoRE from the initial input of a code repository to the final output of a complete repair.

5 Results

5.1 RQ1. Buggy Method and Buggy Lines Localization from Code Repository

5.1.1 Localizing Buggy Method(s). We first use our MethodRecorder tool, which runs the failing tests and records the signatures of the methods that have been executed to narrow down the scope of buggy methods. Then we provide the LLM with the signatures of the executed methods, along with issue content (including the issue description and discussion, which we crawled from each bug's issue URL) and the error stack trace corresponding to the failing test case, as supplied by the Defects4J framework. In this step, we choose not to generate the debugging information. The main reason is that the debugging information can be too long for LLMs to process effectively before the buggy method(s) have been localized. Using the prompt specified in Table 1, we ask the LLM to output the signatures of buggy method(s). To evaluate what kinds of software artifacts can better assist the LLM in localizing buggy methods, we conduct four separate experiments by feeding the LLM with (1) only executed method signatures, (2) executed method signatures and issue content, (3) executed method signatures and error stack trace, and (4) executed method signatures, issue content and error stack trace.

Table 2 shows the effectiveness of buggy method localization in these four experiments. The first row, labeled "—", indicates that only the related (executed) methods are fed to the LLM for buggy method localization. The second, third and fourth rows represent scenarios when, in addition to method localization, issue content, error stack, or both are provided. Following the approach described in LLMAO [41], we consider a buggy method correctly localized if at least one match in the LLM's response corresponds to the buggy method (single-method bugs) or one of the buggy methods (non-single-method bugs) in the ground truth. It is worth mentioning that there exist bug cases where certain software artifacts can not be extracted, and in this case we calculate the ratio by dividing only the number of bug cases with certain software artifacts information. Table 2 shows that, when issue content or error stack is provided, the LLM is able to localize more buggy methods

Table 2. Effectiveness of buggy method localization in Defects4J v2.0 with different software artifacts

	Single Method	Non-Single
_	135/486=27.8%	77/320=24.1%
Issue	207/475=43.6%	127/313=40.6%
Stack	216/486=44.4%	123/320=38.4%
Issue + Stack	234/475=49.3%	149/313=47.6%

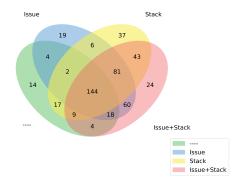


Fig. 7. Overlap of buggy method localization in Defects4J v2.0

Table 3. Method-level FL comparision with baselines on Defects4J v1.2

Projects #Bugs		DEVLoRe		AgentFL [29]		GRACE [25]			FLUCCS [33]				
1 Tojects	#Dugs	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5
Chart	26	10	18	19	16	18	19	14	20	22	15	16	19
Lang	65	47	52	53	44	45	45	42	54	57	40	53	55
Math	106	71	83	84	49	60	61	61	78	89	48	77	83
Time	27	11	14	15	11	13	13	11	14	19	8	15	18
Mockito	38	16	21	22	13	14	14	17	24	26	7	19	22
Closure	133	19	32	35	24	33	35	47	70	81	42	66	77
Total	395	174	220	228	157	183	187	192	260	294	160	249	271
w.o Closure	262	155	188	193	133	150	152	145	190	213	118	183	194

than nothing is provided. When issue content is provided, 15.8% more single buggy methods and 16.5% more non-single buggy methods are correctly localized. When error stack is included, 16.6% more single buggy methods and 14.3% more non-single buggy methods are correctly localized. Furthermore, combining issue content and error stack achieves the best performance, which results in 49.3% in localizing single buggy methods and 47.6% of non-single buggy methods.

Figure 7 shows the overlap of results of buggy method localization with different artifacts in Defects4J v2.0. We can see that when provided with issue content, the LLM can localize 19 extra bugs correctly. Similarly, when provided with error stack, the LLM can localize 37 extra bugs correctly. This finding verifies our assumption in Section 2 that different kinds of software artifacts can complement each other in buggy method localization.

As existing SOTA fault localization approaches use Defects4J v1.2 dataset and adopt Top-n (the buggy method is in the top n list of the LLM's response) to evaluate the performance of fault localization, we compare the performance of DEVLoRe with the SOTA approaches in the same dataset and adopt Top-n approach. Table 3 shows the comparison results where AgentFL [29] is LLM-based, GRACE [25] and FLUCCS [33] are learning-based. As we can see, DEVLoRe outperforms another LLM-based approach AgentFL [29] in Top-1, Top-3 and Top-5 metrics. As discussed in [29], the Closure project contains many bugs with similar code, making it more suitable for learning-based approaches. This characteristic allows GRACE [25] and FLUCCS [33] to achieve an overall high localization rate. The last row shows that, without the Closure project, DEVLoRe outperforms other approaches in Top-1 and achieves similar performance in Top-3 and Top-5, demonstrating its ability to precisely localize buggy methods.

Table 4. Effectiveness of buggy line localization with different software artifacts

	Match	Single Method	Non-Single
	Exact	40.3%	4.3%
_	Range-3	68.1%	8.1%
	Range-5	72.8%	9.8%
	Exact	49.5%	5.2%
Issue	Range-3	74.4%	9.2%
	Range-5	77.5%	11.0%
	Exact	45.0%	4.9%
Stack	Range-3	66.9%	9.0%
	Range-5	71.0%	10.1%
	Exact	39.9%	3.5%
Debug	Range-3	65.4%	6.4%
	Range-5	70.3%	8.1%
	Exact	44.8%	5.2%
Stack + Debug	Range-3	66.9%	9.0%
	Range-5	70.3%	10.7%
	Exact	50.1%	5.2%
Issue + Debug	Range-3	70.3%	8.4%
	Range-5	73.2%	9.8%
	Exact	49.1%	5.5%
Issue + Stack	Range-3	73.2%	10.1%
	Range-5	77.5%	11.8%
	Exact	50.3%	5.2%
Issue + Stack + Debug	Range-3	71.4%	11.3%
	Range-5	73.8%	12.7%
	Exact	68.9%	9.0%
Union	Range-3	83.8%	16.5%
	Range-5	85.7%	19.1%

5.1.2 Localizing Buggy Line(s). To evaluate LLMs' ability in localizing buggy lines, we provide the LLM with the buggy method(s) baseline (i.e., the ground truth which can be extracted from human developer's correct patch) and ask the LLM to predict which lines are buggy. If the line number generated by the LLM exactly matches the first line added in the human developer's correct patch in Defects4J, we consider it an exact match. If the predicted line number falls within a range of n lines from the first added line in the human developer's patch (for example, if the correct patch starts at Line 368, and the LLM's prediction falls within the range from 368 - n to 368 + n), we classify it as a Range-n match.

Table 4 shows the effectiveness of buggy line localization with different artifacts. We can see that when the LLM is fed with issue content, it performs the best among all single software artifacts, achieving 49.5% in exact matches, 74.4% in Range-3, and 77.5% in Range-5 for single methods. When adding more software artifacts, the performance can be further increased. For example, when adding issue content to error stack, the exact, Range-3, and Range-5 buggy line match for single methods can be increased from 45.0% to 49.1%, 66.9% to 73.2%, and 71.0% to 77.5%. The last row indicates that, when combining all artifacts, 68.9% of buggy lines in single-method bugs and 9% of buggy lines in non-single-method bugs are correctly localized. Similar to Figure 7, Figure 8 demonstrates

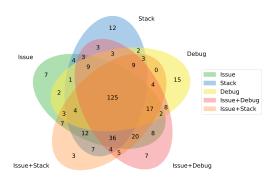


Fig. 8. Buggy line localization exact match with different artifacts in Defects4J v2.0

that, while most of the correctly localized buggy lines across different artifacts overlap, different combinations of artifacts each have their own advantages (since only five areas can be displayed in the overlap figure, we randomly select five categories for visualization). For example, using only the issue content can correctly localize 7 more buggy lines, while using only the debug information can localize 15 more.

Answer to RQ1: Among all three types of single software artifacts, issue content is the most effective in assisting the LLM with bug localization. Furthermore, different types of software artifacts complement each other in localizing buggy methods and buggy lines. This aligns with how human developers typically approach bug localization, as they often rely on information from multiple sources to identify bugs.

5.2 RQ2. Program Repair Based on Provided Method-level Fault Localization

Table 5. Plausible patch with provided method-level fault localization in Defects4J v2.0

	Single Method	Non Single
_	46/489=9.4%	4/345=1.2%
Issue	127/478=26.6%	14/337=4.2%
Stack	132/489=27.0%	13/345=3.8%
Debug	114/468=24.4%	11/326=3.4%
Issue+Stack	159/464=34.3%	21/328=6.4%
Issue+Debug	191/457=41.8%	24/318=7.6%
Stack+Debug	139/462=30.1%	14/321=4.4%
Issue+Stack+Debug	197/457=43.1%	30/318=9.4%
Union	274/489=56.0%	50/346=14.5%

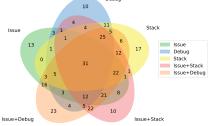


Fig. 9. Overlap of program repair with different artifacts in Defects4J v2.0

To evaluate the effectiveness of program repair, we first extract buggy methods from human's correct patches and used these buggy methods as the baseline. Then we use DebugRecorder to get

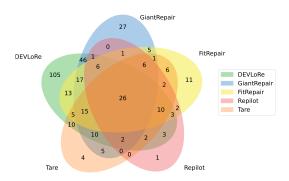


Fig. 10. A Comparison between DEVLoRE with Other State-of-the-art SOTA Program Repair Approaches

the debug information in these buggy methods. Finally, we provide the LLM with different artifacts and check how many plausible patches (passing all tests) it could generate.

Table 5 shows that, error stack achieves the best performance in generating plausible patches for single-method bugs (27.0%) among all three software artifacts, while issue content performs the best for non-single buggy methods. If feeding the LLM with all three types of software artifacts, we can get the best performance, which is 43.1% and 9.4% in patching single-method bugs and non-single-method bugs, respectively. Figure 9 shows that different artifacts complement each other in program repair. Error stack can help the LLM fix 17 extra bugs while issue content can fix 13 extra bugs. Combining issue content with debug info can assist the LLM in fixing 23 extra bugs.

Table 6 compares DEVLoRE with other SOTA program repair approaches. As these approaches focus on repairing single-method bugs, we compare our tool with these approaches on the same dataset. As shown in Table 6, 142 out of 259 bugs in Defects4J v1.2 and 132 out of 230 bugs in Defects4J v2.0 can be repaired by DEVLoRE, outperforming the other five approaches. In total, DEVLoRE can fix 274 buggy methods, which is 60.2% more than GiantRepair [20]. Furthermore, Figure 10 shows that 105 bugs can only be fixed by DEVLoRE. This result demonstrates DEVLoRE's ability in fixing more and extra bugs by leveraging different types of software artifacts.

Table 6. Plausible patch with provided method-level fault localization in single buggy method projects.

Project	#Bugs	DEVLoRe	GaintRepair [20]	FitRepair [36]	Repilot [34]	Tare [47]	GAMMA [44]
Chart	16	12	8	8	6	11	9
Closure	95	37	32	29	21	22	20
Lang	42	32	14	17	15	13	10
Math	74	53	26	23	20	20	19
Time	16	3	1	3	2	3	1
Mockito	16	10	6	4	0	2	2
Defects4J v1.2	259	142	87	85	64	71	61
Defects4J v2.0	230	132	84	44	47	37	39
Total	489	274	171	129	111	108	100

Answer to RQ2: Among all three types of single software artifacts, error stack is the most effective in assisting the LLM with program repair with provided buggy method localization. Different combination of software artifacts complement each other in generating plausible patches and combining all three artifacts can achieve 56.0% of program repair, outperforming the SOTA approaches in program repair.

5.3 RQ3. End-to-end Performance from Code Repository to Program Repair

Currently, most of LLM-based approaches focus on either fault localization or program repair based on the already-localized buggy methods. We investigated RQ1 and RQ2 to compare with these approaches. In contrast, our DEVLoRE fully relies on LLMs for both fault localization and program repair. Furthermore, from fault localization to program repair, we not only provide the LLM with the localized buggy method signatures and lines, but also supply it with the complete method body. When the LLM is asked for the program repair task, it may reconsider and repair other parts of the code beyond the specified buggy location. Therefore, the overall end-to-end performance from the buggy code without localization to a complete program repair cannot be simply calculated by multiplying the fault localization rate by the program repair rate. Conducting a comprehensive end-to-end assessment can offer a better understanding of DEVLoRe's performance throughout the entire fault localization and program repair process.

Table 7 shows DEVLORE's end-to-end performance with different software artifacts. When using issue content alone, the LLM can repair 21.8% of bugs, which is the highest among all single software artifacts. This result is also consistent with Agentless [35], one of the best approaches in SWE-bench lite (with Python projects) [2], which relies solely on issue descriptions to assist LLMs in resolving GitHub issues.

Furthermore, by combining issue content with debugging information and the error stack, the LLM can fix an additional 1.2% and 4.0% of single-method bugs, respectively. When all three software artifacts are combined, the LLM can repair 28.0% of single-method bugs and 11.2% of non-single-method bugs. Furthermore, the combination of different artifacts achieves an end-to-end 39.7% fix rate for single-method bugs and 17.1% for non-single-method bugs. Figure 11 shows similar observations as other two RQs, that different combinations of software artifacts can complement each other in the end-to-end process of bug localization and program repair.

Table 7. End-to-end plausible patch in Defects4J v2.0

	Single Method	Non Single
_	23/488=4.7%	6/346=1.7%
Issue	103/473=21.8%	30/335=9.0%
Debug	38/413=9.2%	3/295=1.0%
Stack	83/488=17.0%	14/346=4.0%
Issue+Debug	92/400=23.0%	30/283=10.6%
Issue+Stack	123/477=25.8%	31/338=9.2%
Stack+Debug	73/413=17.7%	13/295=4.4%
Issue+Stack+Debug	114/407=28.0%	33/295=11.2%
Union	194/489=39.7%	59/346=17.1%

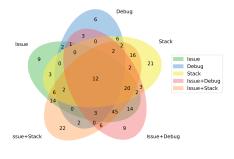


Fig. 11. Overlap of the overall performance in Defects4J v2.0

We divide the time spent on each process by the number of fixed bugs to calculate the average time required to fix a bug. Table 8 shows the average time spent to fix a bug in each process of

Process	Time Spent (On Average)
Extract Related Methods by MethodRecorder	≈ 5s
Localize Buggy Methods by LLM	<1s
Extract Debugging Information by DebugRecorder	$\approx 5s$
Localize Buggy Lines by LLM	<2s
Generate Patches by LLM	<2s
Evaluate Plausible Patches by the Defects4J framework	≈ 300s
Total	≈ 315s

Table 8. Time spent on each process of DEVLoRE

DEVLore using our hardware device (Intel(R) Xeon(R) Platinum 8352V CPU 2.10GHz, 120GB RAM). We observe that the entire localization and repair process takes approximately 15 seconds, with the evaluation of patches being the most time-consuming step, as it involves running all relevant tests. Additionally, since we select the most cost-efficient GPT-40-mini model, the average cost to fix a bug is \$0.057. We calculate this by dividing the total cost of using the LLM when issue content, error stack, and debug information are provided by the number of plausible patches generated, which is 147.

Using our DEVLoRE framework for this end-to-end software maintenance activity, 253 (194 for single buggy methods and 59 for non-single buggy methods) plausible patches can be generated. To the best of our knowledge, DEVLoRE outperforms all current SOTA end-to-end approaches in the Defects4J dataset (the plausible rate of Toggle [8] is 24.2% for 58/240 single-hunk bugs and the plausible rate of FixAgent [19] without Web search engine is 245/835=29.3% at a cost of \$0.364 per bug).

Answer to RQ3: Among the three types of software artifacts, issue content is the most effective in assisting the LLM throughout the end-to-end process, from code repository analysis to fault localization and program repair. Different combinations of software artifacts complement each other, enhancing the overall bug localization and repair process. Moreover, our DEVLoRE framework can effectively fix bugs at a low cost and within a reasonable time frame.

6 Discussion

6.1 Analysis of Results

6.1.1 A simple but efficient framework: Unlike approaches that rely on patch skeletons [20], fix templates [44], or type checking [47], DEVLORE adopts a straightforward framework that allows LLMs to handle both fault localization and program repair with the aid of two lightweight tools we implemented MethodRecorder and DebugRecorder, making it both simple and efficient. The findings of RQ1, RQ2, and RQ3 demonstrate that different software artifacts may lead to different performance in bug localization or program repair. Also, because there is no constraints such as fill-in-the-blank templates or code skeletons for generating patches [20, 44], DEVLORE can fix bugs that other tools cannot, especially some non-single method bugs. Our results in Table 7 show that by combining issue content, stack error, and debug information, our approach can fix 11.2% of non-single method bugs. More importantly, combinations of software artifacts can achieve the best performance in all experiments: bug localization, program repair, and the overall streamlined process. One explanation for the strong performance of DEVLORE is that by integrating multiple

sources of information, the noise inherent in any single source can be significantly reduced [27]. This allows the LLM to focus on the consistent and complementary aspects of the different artifacts, enabling it to perform more effective reasoning and making better use of the available data. While it may be argued that the good results of DEVLoRE are due to the use of GPT-4 models, a recent study shows that directly using the GPT-4 model does not improve the fix rate [20]. Therefore, we argue that our proposed framework is the key to achieving the high rate in fault localization and program repair.

6.1.2 A strict input/output prompt design: Table 1 presents the prompts used in our DEVLORE framework. In the GENERAL TASK PROMPT, the LLM is asked to act as a software engineer and conduct the review process, helping the LLM form a clear understanding of the overall task. The INPUT PROMPT includes various types of software artifacts, with clear symbols denoting different hierarchy levels and structures. For example, the {related methods} in the input prompt wraps the class names in ### symbols, and the method signatures are separated by line breaks. Also, the first line of {debugging_info} in the input prompt represents the currently executed method line, and the second line represents the names and values of the variables in current context (e.g., commons.lang3.math.NumberUtils:createNumber:468[hexDigits:8] represents that the code is about to execute Line 468 and the value of the local variable hexDigits is 8). These strict formatting specifications helps the LLM "understand" the structure of the various information from different software artifacts. The EXPECTED OUTPUT PROMPT is very strict in DEVLORE. When localizing buggy methods, we ask the LLM to return a set of buggy method or field locations in the format path.to.ClassA::methodA. For buggy lines, we request the LLM to return a set of buggy line locations in the format path.to.ClassA line:20. During program repair, we employ the well-known SEARCH/REPLACE method, which is commonly used in many state-of-the-art program repair approaches [23, 28, 35, 45]. By enforcing this strict output format, it significantly reduces the likelihood of hallucinations from the LLM. We believe that the clear and strict prompt design in our DEVLoRE framework has helped the LLM achieve strong performance in fault localization and program repair.

DEVLORE Project GaintRepair [20] Tare [47] TBar [22] Chart 7/127/10 7/10 11/14 Closure 10/16 16/33 12/236/10Lang 8/16 12/1912/194/11 Math 28/48 22/4018/34 12/26 Time 0/21/3 2/3 1/2Mockito 6/116/6 2/21/2Total 59/105 64/111 57/95 31/61 P(%) 56.2% 57.7% 60.0% 50.8%

Table 9. A comparison from plausible patch to correct patch in Defects4J v1.2

6.1.3 Plausible patch to correct patch: Most state-of-the-art (SOTA) program repair approaches use plausible patches that pass all unit tests as an important evaluation metric, since generating plausible patches within limited time and resources is crucial for practical applications. To facilitate a better comparison, we also used plausible patches in RQ2 (Which software artifacts can better assist LLMs in generating valid patches when provided with buggy methods?) and RQ3 (What is the

^{*} X/Y denotes X correct patches and Y plausible patches.

overall performance of DEVLoRE in bug localization and program repair?). Some may argue that a high number of plausible patches does not necessarily translate to a high number of semantically correct patches. To address this, we manually inspected the plausible patches, using the approach described in Section 3.4. Table 9 compares the end-to-end repair results on Defects4J v1.2, as all the baselines consistently used this benchmark. Our DEVLoRE framework uses the LLM to localize buggy methods and lines, whereas all three baselines — GiantRepair [20], Tare [47], and TBar [22] — employ the spectrum-based algorithm Ochiai [3], implemented by GZoltar [4], to localize buggy methods. As shown in Table 9, the ratio of correct to plausible patches (56.2%) produced by our DEVLoRE is lower than that from Tare (60.0%) and from GiantRepair (57.7%). This difference is understandable, as GiantRepair enforces an AST patch skeleton and Tare uses a typing-checking mechanism. Notably, 59 out of 105 plausible patches generated by DEVLoRE are evaluated as correct patches, which is higher than the 57 in Tare and 31 in TBar. This finding demonstrates DEVLoRe's ability to generate candidate patches that not only pass all unit tests but are also semantically correct.

6.2 Threats to Validity

The first threat to validity is the accuracy of the two tools we implemented: *MethodRecorder* and *DebugRecorder*. Both tools rely on mature Java agent technology [1]. We randomly selected several projects and manually verified the outputs of both tools. The manual inspection showed that the outputs were accurate. However, we did not inspect all projects, which may pose a threat to the construct validity.

Another potential threat is that the ChatGPT-40-mini model used in this work may have been trained on open-source projects from GitHub, which could overlap with the Defects4J dataset, leading to possible data leakage. To mitigate this, we also randomly selected 100 bugs from another dataset, GrowingBugs [12], and found the plausible fixing rate to be 39%, which may help alleviate this concern. Additionally, the debugging information requires dynamic analysis, which is unlikely to have been used during the model's training. Also as found in [20], directly using the GPT-4 model can not improve the fix rate.

The final threat to the external validity is that our experiments were conducted on Java projects, and the findings may not be generalized to projects written in other programming languages. To address this, we plan to design *MethodRecorder* and *DebugRecorder* on other programming languages and evaluate DEVLoRE on additional datasets across multiple programming languages in our future work.

7 Related Work

7.1 Large Language Models for Fault Localization

Recently, there has been significant interest in using LLMs for fault localization. Toggle incorporated additional contextual information, such as the buggy line number or code review comments, and greatly enhances the accuracy of predicting both the starting and ending buggy tokens [8]. AGENTFL employs a multi-agent system based on ChatGPT and frames the fault localization task as a three-step process: comprehension, navigation, and confirmation. In each step, AGENTFL deploys specialized agents, each with unique expertise, and uses different tools to address specific tasks [29]. CrashTracker conducts static analysis to map each crash to the corresponding exception instance and identify potential buggy candidates. It then utilizes LLMs to enhance the explainability of the localization results [40]. Jiang et al. assessed the performance of recent commercial closed-source general-purpose LLMs, such as ChatGPT 3.5, ERNIE Bot 3.5, and IFlytek Spark 2.0, on line-level fault localization with the provided buggy method [11]. LLMAO fine-tunes LLMs with 350M, 6B, and

16B parameters on small, curated corpora like Defects4J, improving Top-1 fault localization by 2.3%-54.4% and Top-5 results by 14.4%-35.6%, compared to the state-of-the-art machine learning fault localization [41]. AutoFL prompts an LLM to use function calls for navigating a repository, enabling effective fault localization in large codebases while overcoming the LLM context length limit. It also generates an explanation of the bug and suggests a fault location [14]. LLM4FL combines traditional spectrum-based fault localization with prompt chaining to divide large coverage data into manageable groups. By employing multiple LLM agents, it navigates the codebase more effectively to localize faults [30]. Like LLM4FL and AutoFL which combines transitional fault localization tools, DEVLoRe incorporates a static code analysis tool, *MethodRecorder*, to identify relevant buggy methods in a lightweight manner by invoking failing tests. However, unlike the above approaches, we provide the LLM with a combination of software artifacts, including issue content, error stack traces, and debugging information, which are commonly used by human developers for fault localization. This enables DEVLoRe to outperform state-of-the-art fault localization methods with greater efficiency and lower cost.

7.2 Large Language Models for Program Repair

Recent studies have extensively explored the use of LLMs for program repair. ChatRepair initially provides the LLM with relevant test failure information and then learns from both the failures and successes of previous patching attempts for the same bug, enhancing its ability for more effective APR [39]. GAMMA converts various fix templates into mask patterns and leverages a pre-trained language model to predict the correct code for the masked portions, treating APR as a fill-in-the-blank task [44]. Repilot generates a candidate patch by combining LLM suggestions with a Completion Engine, removing infeasible tokens and filling in gaps proactively [34]. FitRepair integrates the plastic surgery hypothesis into LLM-based APR, combining the direct use of LLMs with two domain-specific fine-tuning strategies and one prompting strategy to enhance its repair capabilities. It can directly generate the correct code in context, effectively "filling in the blanks" of missing code lines or hunks [36]. Tare incorporated type checking into neural program repair model and can successfully repair 62 and 32 bugs from Defects4J v1.2 and Defects4J v2.0 [47]. GiantRepair creates patch skeletons from LLM-generated patches to narrow the patch space, then generates context-aware, high-quality patches by instantiating these skeletons for specific programs [20]. FixAgent unifies debugging through multi-agent collaboration and achieves strong performance with a three-layer hierarchical structure, where the final layer involves the use of a Web search engine [19]. MORepair fine-tunes LLMs for program repair by adapting both to the syntactic nuances of code transformation and the underlying logic of code changes, enabling the generation of high-quality patches [43]. To leverage LLMs' capabilities and augmented information, CREF is a semi-automatic repair framework for programming tutors, highlighting the potential for enhancing LLMs' repair capabilities through tutor interactions and historical conversations [42]. RepairLLaMA is an innovative program repair method that finds optimal code representations for APR using fine-tuned models, and introduces a state-of-the-art parameter-efficient fine-tuning technique (PEFT) for program repair [31]. Our approach differs from the aforementioned approaches in two key ways. First, we do not rely on fill-in-the-blank templates or skeletons for generating patches, which helps avoid patch overfitting problems in program repair [6, 16, 32]. Second, our fault localization step allows for the identification of multiple buggy methods, which are then fed into the LLM's repair process along with various software artifacts. This provides flexibility to address not only single-method bugs but also bugs spanning across different methods. As seen in Table 7, the end-to-end process from localization to repair can generate plausible patches for 17.1% of non-single-method bugs, while most existing program repair approaches primarily target single-method bugs.

8 Conclusions and Future Work

This paper presents an LLM-based framework, DEVLORE, for streamlining fault localization and program repair. By mimicking human developers in addressing bug problems and integrating three different software artifacts, DEVLORE demonstrates strong performance in both fault localization and program repair, outperforming current state-of-the-art approaches in terms of bug fixing rate, time, and cost. In addition, unlike more rigid approaches that ask LLMs to fill in the blank within a single buggy method or use the fix template, DEVLORE feed the LLM with only different software artifacts and there are no constraints on the DEVLORE framework regarding how it repairs bugs, DEVLORE has shown significant potential in handling bugs that span across multiple methods.

Our future work will focus on the following directions: first, expanding DEVLoRE to support additional programming languages, such as Python and C/C++, and evaluating its effectiveness on projects written in these languages; and second, testing the DEVLoRE framework with a broader range of software artifacts, such as commit history, code review comments, and user documentation, to assess how these additional data sources can further enhance the bug localization and program repair process.

Data Availability

The source code and experimental results of this work for replication are available at https://github.com/XYZboom/DEVLoRe.

Acknowledgments

This work has been partially supported by the National Natural Science Foundation of China (NSFC) with Grant No. 62172311.

References

- [1] [n. d.]. Instrumentation documents in Java SE 8. https://docs.oracle.com/javase/8/docs/api/java/lang/instrument/ Instrumentation.html
- [2] 2024. SWE-bench Lite. https://www.swebench.com/lite.html/.
- [3] Rui Abreu, Peter Zoeteweij, and Arjan JC Van Gemund. 2007. On the accuracy of spectrum-based fault localization. In Testing: Academic and industrial conference practice and research techniques-MUTATION (TAICPART-MUTATION 2007). IEEE, 89–98.
- [4] José Campos, André Riboira, Alexandre Perez, and Rui Abreu. 2012. Gzoltar: an eclipse plug-in for testing and debugging. In *Proceedings of the 27th IEEE/ACM international conference on automated software engineering*. 378–381.
- [5] Zimin Chen, Steve Kommrusch, Michele Tufano, Louis-Noël Pouchet, Denys Poshyvanyk, and Martin Monperrus. 2019. Sequencer: Sequence-to-sequence learning for end-to-end program repair. IEEE Transactions on Software Engineering 47, 9 (2019), 1943–1959.
- [6] Zhiwei Fei, Jidong Ge, Chuanyi Li, Tianqi Wang, Yuning Li, Haodong Zhang, LiGuo Huang, and Bin Luo. 2024. Patch Correctness Assessment: A Survey. ACM Transactions on Software Engineering and Methodology (2024).
- [7] 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). ACM, 19–30.
- [8] Soneya Binta Hossain, Nan Jiang, Qiang Zhou, Xiaopeng Li, Wen-Hao Chiang, Yingjun Lyu, Hoan Nguyen, and Omer Tripp. 2024. A Deep Dive into Large Language Models for Automated Bug Localization and Repair. Proceedings of the ACM on Software Engineering 1, FSE (2024), 1471–1493.
- [9] Nan Jiang, Kevin Liu, Thibaud Lutellier, and Lin Tan. 2023. Impact of code language models on automated program repair. In *Proceedings of the 45th IEEE/ACM International Conference on Software Engineering (ICSE)*. IEEE, 1430–1442.
- [10] 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). IEEE, 1161–1173.
- [11] Shengbei Jiang, Jiabao Zhang, Wei Chen, Bo Wang, Jianyi Zhou, and Jie Zhang. 2024. Evaluating Fault Localization and Program Repair Capabilities of Existing Closed-Source General-Purpose LLMs. In Proceedings of the 1st International Workshop on Large Language Models for Code (LLM4Code). ACM, 75–78.
- [12] Yanjie Jiang, Hui Liu, Xiaoqing Luo, Zhihao Zhu, Xiaye Chi, Nan Niu, Yuxia Zhang, Yamin Hu, Pan Bian, and Lu Zhang. 2022. Bugbuilder: An automated approach to building bug repository. *IEEE Transactions on Software Engineering* 49, 4

- (2022), 1443-1463.
- [13] 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)*. ACM, 437–440.
- [14] Sungmin Kang, Gabin An, and Shin Yoo. 2024. A Quantitative and Qualitative Evaluation of LLM-Based Explainable Fault Localization. Proceedings of the ACM on Software Engineering 1, FSE (2024), 1424–1446.
- [15] Xuan Bach D. Le, David Lo, and Claire Le Goues. 2016. History driven program repair. In *Proceedings of the 23rd IEEE International Conference on Software Analysis, Evolution, and Reengineering (SANER)*. IEEE, 213–224.
- [16] Xuan-Bach D Le, Ferdian Thung, David Lo, and Claire Le Goues. 2018. Overfitting in semantics-based automated program repair. In *Proceedings of the 40th international conference on software engineering*. 163–163.
- [17] Claire Le Goues, ThanhVu Nguyen, Stephanie Forrest, and Westley Weimer. 2011. Genprog: A generic method for automatic software repair. IEEE Transactions on Software Engineering 38, 1 (2011), 54–72.
- [18] Claire Le Goues, Michael Pradel, and Abhik Roychoudhury. 2019. Automated program repair. *Commun. ACM* 62, 12 (2019), 56–65.
- [19] Cheryl Lee, Chunqiu Steven Xia, Longji Yang, Jen-tse Huang, Zhouruixin Zhu, Lingming Zhang, and Michael R Lyu. 2024. A unified debugging approach via llm-based multi-agent synergy. arXiv preprint arXiv:2404.17153 (2024).
- [20] Fengjie Li, Jiajun Jiang, Jiajun Sun, and Hongyu Zhang. 2024. GiantRepair: Hybrid Automated Program Repair by Combining Large Language Models and Program Analysis. arXiv preprint arXiv:2406.00992 (2024).
- [21] 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)*. ACM, 602–614.
- [22] 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)*. ACM, 31–42.
- [23] Yizhou Liu, Pengfei Gao, Xinchen Wang, Jie Liu, Yexuan Shi, Zhao Zhang, and Chao Peng. 2024. MarsCode Agent: AI-native Automated Bug Fixing. arXiv preprint arXiv:2409.00899 (2024).
- [24] Fan Long and Martin Rinard. 2015. Staged program repair with condition synthesis. In *Proceedings of the 10th Joint Meeting on Foundations of Software Engineering (FSE)*. ACM, 166–178.
- [25] Yiling Lou, Qihao Zhu, Jinhao Dong, Xia Li, Zeyu Sun, Dan Hao, Lu Zhang, and Lingming Zhang. 2021. Boosting coverage-based fault localization via graph-based representation learning. In Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering. 664–676.
- [26] 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). ACM, 691–701.
- [27] Xiangxin Meng, Zexiong Ma, Pengfei Gao, and Chao Peng. 2024. An Empirical Study on LLM-based Agents for Automated Bug Fixing. arXiv preprint arXiv:2411.10213 (2024).
- [28] Siru Ouyang, Wenhao Yu, Kaixin Ma, Zilin Xiao, Zhihan Zhang, Mengzhao Jia, Jiawei Han, Hongming Zhang, and Dong Yu. 2024. RepoGraph: Enhancing AI Software Engineering with Repository-level Code Graph. arXiv preprint arXiv:2410.14684 (2024).
- [29] Yihao Qin, Shangwen Wang, Yiling Lou, Jinhao Dong, Kaixin Wang, Xiaoling Li, and Xiaoguang Mao. 2024. AgentFL: Scaling LLM-based Fault Localization to Project-Level Context. arXiv preprint arXiv:2403.16362 (2024).
- [30] Md Nakhla Rafi, Dong Jae Kim, Tse-Hsun Chen, and Shaowei Wang. 2024. Enhancing Fault Localization Through Ordered Code Analysis with LLM Agents and Self-Reflection. arXiv preprint arXiv:2409.13642 (2024).
- [31] André Silva, Sen Fang, and Martin Monperrus. 2023. Repairllama: Efficient representations and fine-tuned adapters for program repair. arXiv preprint arXiv:2312.15698 (2023).
- [32] Edward K Smith, Earl T Barr, Claire Le Goues, and Yuriy Brun. 2015. Is the cure worse than the disease? overfitting in automated program repair. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*. 532–543.
- [33] Jeongju Sohn and Shin Yoo. 2017. Fluccs: Using code and change metrics to improve fault localization. In *Proceedings* of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA). ACM, 273–283.
- [34] Yuxiang Wei, Chunqiu Steven Xia, and Lingming Zhang. 2023. Copiloting the Copilots: Fusing Large Language Models with Completion Engines for Automated Program Repair. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESC/FSE)*. ACM, 172–184.
- [35] Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2024. Agentless: Demystifying LLM-based Software Engineering Agents. arXiv preprint arXiv:2407.01489 (2024).
- [36] Chunqiu Steven Xia, Yifeng Ding, and Lingming Zhang. 2023. The Plastic Surgery Hypothesis in the Era of Large Language Models. In *Proceedings of the 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 522–534.

[37] Chunqiu Steven Xia, Matteo Paltenghi, Jia Le Tian, Michael Pradel, and Lingming Zhang. 2024. Fuzz4All: Universal Fuzzing with Large Language Models. arXiv preprint arXiv:2308.04748 (2024).

- [38] 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 IEEE/ACM International Conference on Software Engineering (ICSE)*. IEEE, 1482–1494.
- [39] 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).
- [40] Jiwei Yan, Jinhao Huang, Chunrong Fang, Jun Yan, and Jian Zhang. 2024. Better Debugging: Combining Static Analysis and LLMs for Explainable Crashing Fault Localization. arXiv preprint arXiv:2408.12070 (2024).
- [41] Aidan Z. H. Yang, Claire Le Goues, Ruben Martins, and Vincent Hellendoorn. 2024. Large Language Models for Test-Free Fault Localization. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering (ICSE)*. ACM, 1–12.
- [42] Boyang Yang, Haoye Tian, Weiguo Pian, Haoran Yu, Haitao Wang, Jacques Klein, Tegawendé F Bissyandé, and Shunfu Jin. 2024. Cref: An LLM-based conversational software repair framework for programming tutors. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*. ACM, 882–894.
- [43] Boyang Yang, Haoye Tian, Jiadong Ren, Hongyu Zhang, Jacques Klein, Tegawendé F Bissyandé, Claire Le Goues, and Shunfu Jin. 2024. Multi-Objective Fine-Tuning for Enhanced Program Repair with LLMs. arXiv preprint arXiv:2404.12636 (2024).
- [44] Quanjun Zhang, Chunrong Fang, Tongke Zhang, Bowen Yu, Weisong Sun, and Zhenyu Chen. 2023. Gamma: Revisiting template-based automated program repair via mask prediction. In *Proceedings of the 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 535–547.
- [45] Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. 2024. Autocoderover: Autonomous program improvement. In Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis. 1592–1604.
- [46] Qihao Zhu, Zeyu Sun, Yuan-an Xiao, Wenjie Zhang, Kang Yuan, Yingfei Xiong, and Lu Zhang. 2021. A syntax-guided edit decoder for neural program repair. In Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE). ACM, 341–353.
- [47] Qihao Zhu, Zeyu Sun, Wenjie Zhang, Yingfei Xiong, and Lu Zhang. 2023. Tare: Type-aware neural program repair. In Proceedings of the 45th IEEE/ACM International Conference on Software Engineering (ICSE). IEEE, 1443–1455.