

# LLMs as Continuous Learners : Improving the Reproduction of Defective Code in Software Issues

YALAN LIN<sup>\*</sup>, Shanghai Jiao Tong University, China

YINGWEI MA, RONGYU CAO, BINHUA LI, FEI HUANG, Tongyi Lab, Alibaba Group, China

XIAODONG GU<sup>†</sup>, Shanghai Jiao Tong University, China

YONGBIN LI<sup>†</sup>, Tongyi Lab, Alibaba Group, China

Reproducing buggy code is the first and crucially important step in issue resolving, as it aids in identifying the underlying problems and validating that generated patches resolve the problem. While numerous approaches have been proposed for this task, they primarily address common, widespread errors and struggle to adapt to unique, evolving errors specific to individual code repositories. To fill this gap, we propose EvoCoder, a multi-agent continuous learning framework for issue code reproduction. EvoCoder adopts a reflection mechanism that allows the LLM to continuously learn from previously resolved problems and dynamically refine its strategies to new emerging challenges. To prevent experience bloating, EvoCoder introduces a novel hierarchical experience pool that enables the model to adaptively update common and repo-specific experiences. Our experimental results show a 20% improvement in issue reproduction rates over existing SOTA methods. Furthermore, integrating our reproduction mechanism significantly boosts the overall accuracy of the existing issue-resolving pipeline.

CCS Concepts: • **Software and its engineering** → **Software testing and debugging**.

Additional Key Words and Phrases: Issue Resolving, Issue Reproduction, Continuous Learning, Large Language Models, Software Engineering Agent

## 1 Introduction

Issue resolving is a fundamental aspect of software development and maintenance, critical to preserving the quality and stability of software systems [26, 41, 49, 53]. Throughout a project’s life-cycle, various issues inevitably arise. Automating the resolving of these issues not only accelerates error identification and correction but also boosts development efficiency, playing a pivotal role in sustaining the long-term health of the project [16, 28, 40, 50]. This process is especially vital in the early stages of development, where test cases are often incomplete, requiring continuous refinement as new bugs are reported by users [17].

In issue resolving, the initial and critical step is issue reproduction which involves automatically generating code to reproduce the reported problem based on the issue description. Specifically, given an issue from a code repository, which contains key steps to reproduce the problem (as shown in Figure 1), the goal is to generate executable code to replicate the issue. The output of this generated code should match the “Actual Result” specified in the issue while applying an appropriate patch should modify the code’s output to reflect the “Expected Result”. Successful issue reproduction not only accelerates the process of problem localization and resolving but also strengthens quality assurance processes in continuous integration and delivery, thereby improving the robustness and reliability of software systems.

Previous studies have extensively explored the role of reproducing code [8, 44, 50]. Some research integrates this process into code intelligent agents (e.g., CodeR [8] and SWE-Agent [44]). These methodologies design computer interfaces for viewing, searching, and editing files, thereby enabling

---

<sup>\*</sup> Work during Yalan’s internship at Tongyi Lab@Alibaba group.

<sup>†</sup> Corresponding Authors.

Authors’ Contact Information: Yalan Lin<sup>\*</sup>, Shanghai Jiao Tong University, China; Yingwei Ma, Rongyu Cao, Binhua Li, Fei Huang, mayingwei.myw@alibaba-inc.com, Tongyi Lab, Alibaba Group, China; Xiaodong Gu<sup>†</sup>, xiaodong.gu@sjtu.edu.cn, Shanghai Jiao Tong University, China; Yongbin Li<sup>†</sup>, Tongyi Lab, Alibaba Group, China.

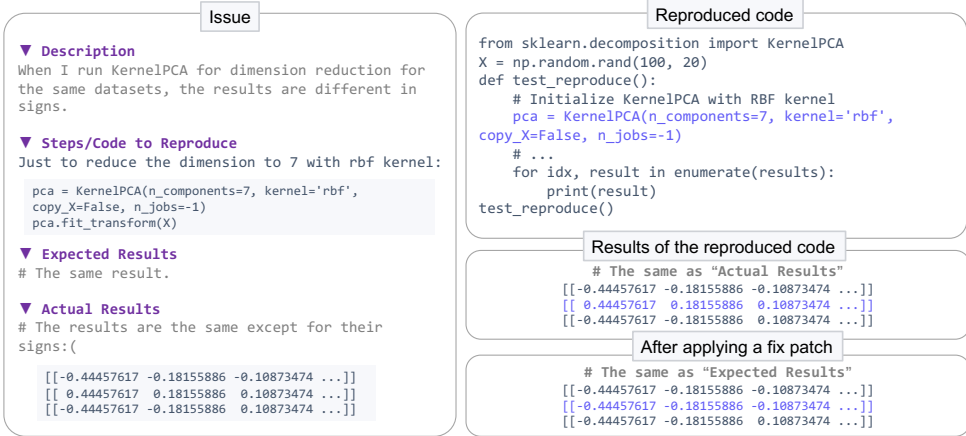


Fig. 1. An example of issue reproduction for the KernelPCA library. The *Expected Results* section specifies the ground truth results for the original problem. The *Actual Results* section specifies the erroneous results that need to be reproduced based on the given steps. As seen in the upper right, the reproduced code successfully produces the “Actual result” in the issue (the three vectors are different in sign). To resolve this issue, a patch is needed to fix the bug. As seen in the lower right, applying the patch aligns the result with the “Expected result”, ensuring that all three vectors are now correctly uniform in sign.

LMs to resolve issues by chatting with computers as humans do. In contrast, agent-less approaches such as AutoCodeRover [50] and Agentless [40] generate test cases from issues, and use test cases for filtering and debugging the reproduced code.

Despite these advancements, current solutions primarily address common, widespread errors and struggle to adapt to unique, evolving errors specific to individual code repositories. For example, some issues depend on project-specific formatting conventions to accurately reproduce. This might involve a pipeline where dependency libraries are introduced first, followed by defining the reproduction function within a *test\_reproduce* function, and subsequently calling this function. Additionally, some issues rely on human-defined rules that LLMs struggle to interpret and apply effectively. These project-specific conventions are typically unique to each repository and, as such, are not easily transferable across different projects. A more comprehensive illustration of the motivation will be presented in Section 3.

In this research, we propose EvoCoder, a novel continual learning framework for issue code reproduction. EvoCoder introduces a hierarchical reflection architecture where an actor LLM reproduces the issue code and stores the action trajectory in the memory. A Reflection LLM distills experiences from its trajectory. To better maintain and utilize the extracted experience, we design a hierarchical experience pool: the higher layer stores general experiences while the lower layer corresponds to repository-specific experiences. To keep the experiences up-to-date and continuously refined, we define five actions including ADD, REMOVE, MERGE, ENDORSE, and MODIFY, for the reflection LM to manipulate the hierarchy of experiences. Our approach enables the model to continuously learn and optimize as it encounters new problems, without incurring additional computational costs associated with fine-tuning. In addition, it dynamically updates and refines the model’s existing knowledge base, allowing it to develop expertise within a specific domain.

We evaluate the effectiveness of the EvoCoder on SWE-bench [16] and compare it with the state-of-the-art issue code reproduction methods. Experimental results reveal that EvoCoder can improve the reproduction rate by 20%. In addition, the generated reproduction code successfully assists in the model’s debugging process, thereby enhancing the accuracy of issue resolving.

The main contributions of our paper are threefold:

- An in-depth analysis of the root causes behind the failure of previous models on issue reproduction.
- We propose a novel issue reproduction method based on continual learning and reflection LM, resulting in a 20% increase in success rate.
- Extensive evaluation of the proposed method on state-of-the-art issue resolving benchmarks.

## 2 Background

### 2.1 Issue Resolving

Issue resolving is a critical component of software development and maintenance, encompassing the identification, diagnosis, and resolving of software defects or issues reported by users or developers. In this task, based on the provided problem descriptions and snapshots of the codebase, the model needs to autonomously identify the specific locations that need modification [14, 29, 38, 50]. This can be achieved through commands to search for certain files or by analyzing execution paths using existing passing and failing test cases to pinpoint the most likely points of failure. Once these locations are identified, the model generates the corresponding code snippets to address the issues [25, 39, 46, 52]. During the generation process, the model can also perform debugging and iterative revisions on its own [10, 12, 31, 35, 43]. For instance, syntax checkers can identify syntactical errors, and test cases or code snippets that reproduce the issue can be executed to verify whether the generated patches resolve the original errors. Any error messages or feedback from these processes can be fed back into the model to further refine and improve the generated code. In the final testing phase, the generated code patches must pass both the unit tests associated with the current modifications—extracted from the pull requests submitted by programmers—and the existing unit tests. This ensures that the new changes do not adversely affect previously implemented functionalities, the generated code patches must pass both the unit tests associated with the current modifications—extracted from the pull requests submitted by programmers—and the existing unit tests. This ensures that the new changes do not adversely affect previously implemented functionalities.

### 2.2 Issue Code Reproduction

Issue Reproduction involves the creation of executable code designed to reproduce the issues reported by users. Historically, this challenge has been framed as a single-step code generation task, wherein the model processes the issue description provided by the user and outputs the corresponding code snippet intended to recreate the reported problem [17]. However, with advancements in LLMs, there has been a shift towards models capable of engaging in multi-turn dialogues. This evolution allows the models to not only generate initial code but also to interact with a simulated environment for debugging purposes. Prominent instances of this advanced approach can be found in systems like SWE-Agent [44] and CodeR [8], both of which utilize agent-based methodologies. In these frameworks, the interaction begins with the model receiving a system prompt that outlines the range of actions it can undertake, including edit, search, and other relevant tasks. Subsequently, the user inputs the specifics of the issue they wish to address. Operating in a ReAct paradigm [45], the model proceeds to analyze the current state, contemplate potential steps, and execute the most suitable action. This typically involves generating the necessary files, modifying them to integrate

the required features, and systematically debugging the program until it faithfully replicates the 'actual result' specified in the original issue report, or halting if multiple attempts fail to yield further progress.

In our work, we continue to build upon this agent-based paradigm, exploring how models can achieve sustained improvement in the process of reproduce a wide variety of issues, thereby refining its ability to accurately and efficiently reproduce issues over time.

### 3 Motivation

In this section, we motivate our approach by analyzing the failure cases of CodeR [8], the state-of-the-art method for issue reproduction. We run CodeR on SWE-bench-lite [6] and extract reproduction code generated by CodeR from its resolving trajectories. Upon running the generated issue code, we manually inspect the execution output and assess its correctness based on the following criteria:

- (1) **Completeness:** The reproduced code must contain the core or complete code provided by the users in the original issue.
- (2) **Consistency:** (i) The error messages upon running the code should be consistent with the provided full error messages. If the input issue involves adding a new feature, the output should contain the new feature. (ii) Execution commands involved in the code must match those described in the problem statement exactly.
- (3) **Authenticity:** Error messages must be derived from actual runtime results, instead of mocked outputs.

Ultimately, we collect a total of 84 erroneous codes generated by CodeR. Through manual examination, we identified seven types of errors that can be categorized into two groups:

**Internal errors (58.55%)** refers to intrinsic errors that arise from the reproduction code itself, including 1) *Wrong Reproduction Targets (4.00%)*: The execution results from the reproduced code do not align with the expected error message described in the original issue. This is likely due to the model's difficulty in fully interpreting the natural language within the issue description—such as distinguishing between error outputs and correct outputs. 2) *Wrong Function Call (10.97%)*: the reproduced code calls wrong functions or commands compared to the intended issue code. 3) *Over-mocking (14.63%)*: Mocking is a commonly used technique in software testing that isolates test scripts from external dependencies by replacing them with mocked functions. However, when the behavior of these dependencies is overly predefined, the mocked functions may not accurately reflect the actual behavior of the original code. For instance, using print statements to output error messages instead of triggering real errors can cause tests to pass even if underlying defects exist. This limits the tests' effectiveness in detecting issues, especially when the true behavior of dependencies changes. 4) *Missing Environment Requirements (6.10%)*: The reproduction only contains the core code while missing the environmental setup. 5) *Inaccurate Execution Results (23.17%)*: The reproduced code contains logical errors during the reproduction process, causing inconsistent results to those described in the issue. This is probably due to the restriction of model capability and the task complexity.

**External errors (41.45%)** refers to the errors from external environments or executions, including 1) *Incorrect commands to run the code (14.63%)*: Some issues require specific commands for reproduction, which are usually mentioned in the issue description or pertain to particular usages of the repository. However, the previous methods uniformly used Python commands to run the code, which could result in the inability to reproduce certain problems and 2) *Wrong environment setting (26.82%)*: Given that reproducing these issues often requires specific versions of libraries, particular operating systems, and interactive environments (such as Jupyter), it may happen that

even with correct code and execution instructions, the reported problem cannot be reproduced if the exact environment is not set up.

While the external errors can be easily handled (e.g., by wrapping the reproduced code with a shell script, including the required environmental modifications and commands for execution), our focus shifts to the more complex internal errors—often the most challenging aspect of issue reproduction. Through a comprehensive analysis of errors encountered during code reproduction, we identified two recurring patterns:

| django_django-15388  | pytest-dev_pytest-7220  | sphinx-doc_sphinx-8435  |
|--|---|---|
| <pre>def test_reproduce():     # Simulate running the     development server and saving a     file     print("Running development     server...")     .....     test_reproduce()</pre> | <pre>def test_reproduce():     .....     try:         assert False     except AssertionError:         print(f"Displayed: ...")     test_reproduce()</pre> | <pre>def test_reproduce():     # Simulating the Sphinx     autodoc behavior     print(f"var type:     {type(var).__name__}")     .....     test_reproduce()</pre> |

(a) Example of the common errors across different libraries.

| django_django-10924   | django_django-14580  | django_django-14855  |
|---|--|--|
| <pre>Traceback (most recent call last): ..... ModuleNotFoundError: No module named 'base'</pre> | <pre>Traceback (most recent call last): ..... ModuleNotFoundError: No module named 'app'</pre> | <pre>Traceback (most recent call last): ..... ModuleNotFoundError: No module named 'myapp'</pre> |

(b) Example of recurring errors within the same library.

Fig. 2. Examples of the characteristics of the encountered errors.

**Finding 1:** Some errors are common across libraries, such as cases that output simulated error messages without reproducing the actual issue (Figure 2a). Resolving these requires the model to continuously distill general rules from prior resolving experiences.

**Finding 2:** Some errors are library-specific, recurring within particular libraries. For example, the errors illustrated in Figure 2b, which are frequently occurred, are often caused by the initial setup in the Django library. Resolving these issues requires the LLM to adaptively acquire specialized knowledge for each repository.

These insights are akin to how human programmers gradually become proficient with a library by repeatedly using the library, solving its issues, and accumulating generic knowledge across different libraries. This motivates us with a continuous learning strategy for issue reproduction, where each reproduction attempt accumulates experience, fostering higher levels of automation and broader applicability in issue reproduction.

## 4 Method

### 4.1 Overview

In this paper, we propose a continuous learning pipeline that enables LLM agents to accumulate experience from previously encountered issues. Unlike conventional methods, our approach does not require fine-tuning; instead, it facilitates the model’s ability to continuously update and optimize its stored experiences when addressing new challenges. The overall methodology is structured into three main components: (i) an *Actor LM* which reproduces the issue using instructions and previous experiences; (ii) a *Reflection LM* which extracts experience from the Actor’s reproduction trajectories; and (iii) a hierarchical *Experience Pool* which stores general and repository-specific experiences, enabling the Reflection LM to continuously update and refine its accumulated knowledge.

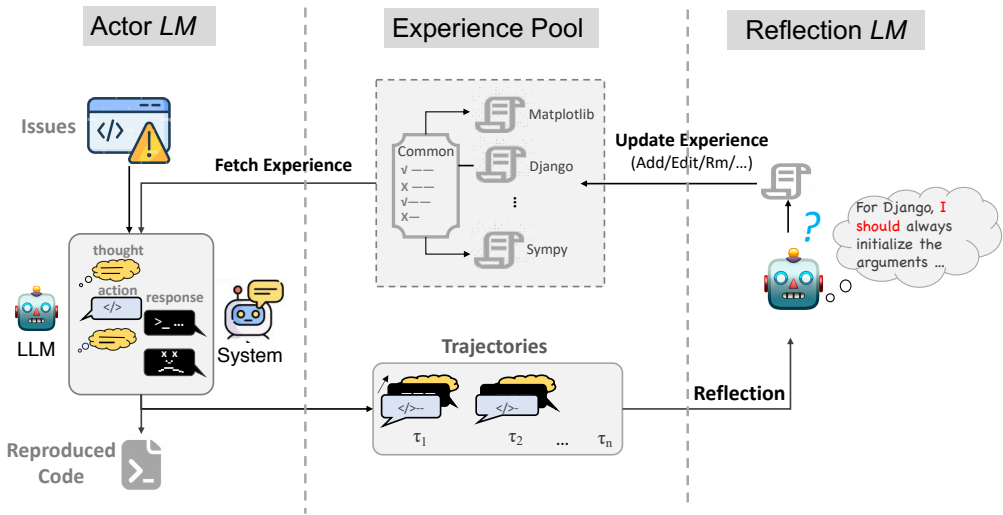


Fig. 3. The overall framework of our approach

### 4.2 Actor: Reproduction Trial

The Actor is an LLM prompted to perform issue reproduction. Adopted from SWE-agent [44] and CodeR [8], we formulate issue reproduction as a multi-turn conversation between agent and computer. In each turn, the Actor is presented with a task instruction, the current system’s output, and experiences from past resolving. It is asked to output a “Thought” about the current reproducing step, followed by an “Action” like searching, viewing, and editing files. Once the reproduction “Action” has been executed, the system responds with the new output, which is fed back to the Actor as input for the next interaction. The dialogue continues, yielding a trajectory of Thought-Action-Response sequences. To mitigate problems caused by improper environment configuration and instruction misunderstanding, we also introduce a standardized process, which asks the model to write instructions to run the code or install related packages in a sh script. The detailed prompt for the Actor is as follows:

## Prompt for Actor

We're currently solving the following issue within our repository. Here's the issue text:  
ISSUE: <issue>

### INSTRUCTIONS:

You are now going to reproduce the provided issue (not solve the issue). Begin your terminal session in the root directory of the repository. To assist you, use any bash commands or the special interface. Make sure to open and read all the files you need, and execute any tests necessary to reproduce the issue. Remember, YOU SHOULD ENTER ONE COMMAND AT A TIME. Always wait for a response after each command before proceeding.

Once you have successfully reproduced the issue and are ready to report it, you can record the steps you took. However, note that you cannot use any interactive session commands (e.g. python, vim) in this environment, but you can run scripts. For example, you can execute a python script with python <script\_name>.py. But DO NOT run './run\_reproduce.sh' or 'bash ./run\_reproduce.sh', use run\_reproduce\_code instead.

Once the reproduction is complete, please output the address of the file containing the reproduction script in the following format.

NOTE ABOUT THE EDIT COMMAND: Indentation really matters! When editing a file, make sure to insert appropriate indentation before each line! You should also check which file you opened in the editor before editing it.

### IMPORTANT TIPS:

1. Always start by trying to replicate the bug that the issues discusses.

If the issue includes code for reproducing the bug, we recommend that you re-implement that in your environment. You MUST create files and execute follow the description, and run it to make sure you can reproduce the bug.

If the bug reproduction script does not print anything when it successfully runs, we recommend adding a print("Script completed successfully, no errors.") command at the end of the file, so that you can be sure that the script indeed ran fine all the way through.

2. When reproducing the code, you should consider all the cases mentioned in the issue.

Before returning, check whether your reproduction of test is complete. The test should be consistent with the issue description. Do not miss the content to be tested. NOTE that the provided code in the issue description of test may be PARTIAL, please generate a complete version of the test code based on the description.

3. If the bug reproduction script requires inputting/reading a specific file, such as buggy-input.png, and you'd like to understand how to input that file, conduct a search in the existing repo code, to see whether someone else has already done that. Do this by running the command: find\_file "buggy-input.png" If that doesn't work, use the linux 'find' command.

4. If you are uncertain about the specific line number when searching for a snippet of code, a class, or a method within a file, prioritize using the grep -nr <code/class/method>command to retrieve the approximate location. Then, after opening the file, use the gotocommand to navigate to the snippet, preventing excessive use of the scroll downcommand. If the grep -nrcommand does not retrieve any relevant content, consider using the scroll downcommand to search for the code after opening the file.

5. During the reproduction process, if you cannot reproduce due to missing packages in the environment, you MUST use commands like pip, apt-get -y, etc. to install the corresponding packages, please assign exact version of packages if you know when you use pip. Please use -quiet to suppress the output when installing packages. DO NOT install <repo> package as it is already installed and fixed. Remember to write these commands into run\_reproduce.sh.

6. You MUST create a file 'run\_reproduce.sh' under the ROOT path of this repo and write the commands into it, including packages you need to install. If you don't need them, simply write commands for code execution into it. Then you can use 'run\_reproduce\_code' to run the code.

7. You must execute and modify the file UNTIL you can reproduce the issue. You can create a reproduce.py script for this purpose. If you are able to run the commands directly within a shell script or create other files follow the issue description, then this script is NOT necessary. The structure of your reproduce.py code should output in the following format:

```
...
import ...
# setup test environment here
```

```

# here is core test code, you MUST use 'test_reproduce' as function name.
def test_reproduce():
<core code>...

# you MUST call 'test_reproduce' here.
test_reproduce()

...

8. After editing the files, use 'run_reproduce_code' to check whether the issue is reproduced. Please use
'open <file_name>' to make sure the modified reproduce file have correct logic and STILL satisfies the
above format requirements.

9. If this issue pertains to the addition of a new feature, your reproduction code should serve to test
whether this functionality has been implemented.

10. You MUST show the output to the console. If you really need to output the result to a file, you MUST
use the 'cat' command to output the result to the console.

11. This is the real environment. Please DO NOT use mock or simulation methods to solve the issue. You
need to reproduce the issue exactly as described in the issue.

12. If you think you have already reproduced the issue, you MUST use 'check' to check whether the issue
is reproduced correctly before you submit the code.

13. Only return "submit" command when the current task is truly completed. If the current task still
cannot be completed after trying many different solutions, please return "fail" command.

13. Here are some experiences summarized from other issues in the same repository. Please refer to these
experiences during generation to avoid making the same mistakes.

<experience>

You MUST take into account all the cases mentioned in the issue and reproduce them.

```

Upon reproducing the entire issue, the dialogue history, known as the resolving trajectory, is stored in a short-term memory. The trajectory contains concrete steps, thoughts, and strategies for resolving issues, providing valuable guidelines for issue resolving. For example, some issues could provide complete reproduction code, including specific usages of the library that the model might not be able to generate on its own. The model undergoes debugging steps during the multi-round dialogues, which can guide it in avoiding initial errors or steering clear of unproductive debugging directions.

### 4.3 Reflection LM: Extracting Experience from Reproduction Trajectories

Having collected the trajectories, we distill them into experiences, guiding further issue reproduction. An *experience* is defined as a rule guiding what the Agent must do to avoid a certain reproduction failure or follow unique coding styles that certain code repositories may have. While previous work [23, 37] proposed using experience pools for knowledge accumulation, their approach faces two critical challenges: (1) experience bloat, where experiences become increasingly verbose and detailed over time, hindering accurate issue pattern matching and experience utilizing, and (2) experience rigidity, where experiences become static and fail to adapt to emerging issue types and repository-specific error patterns.

To enable effective experience refinement, we design an LLM-based reflection mechanism that analyzes task trajectories and golden test patches to continuously update the experience pool. Unlike previous approaches that simply accumulate experiences, our method actively manages experience quality through carefully designed prompts that instruct the reflection LLM to: (1) Analyze Trial. Examines the reproduction trajectory and compares it with the golden test patch to identify



successful patterns or failure causes. This analysis determines whether to: ADD new experiences when discovering novel reproduction patterns; MODIFY existing experiences when current solutions can be improved; REMOVE experiences that consistently lead to failed reproductions. (2) Analyze Rule Applicability. Evaluate whether the experience should be: AGREED UPON and kept if it proves effective across multiple repositories; MERGED with similar experiences to maintain conciseness; Categorized as repository-specific or generally applicable.

For example, when analyzing a successful reproduction trial, the LLM first compares the issue description with the golden test patch to understand the core reproduction pattern. If this pattern represents a novel approach, it triggers an Add operation; if it refines an existing pattern, it initiates a Modify operation. Through these reflection-based operations, the model continuously improves its reproduction capability in two ways: First, by maintaining a curated set of high-quality experiences that accurately capture issue reproduction patterns; Second, by organizing these experiences into appropriate scopes (general or repository-specific) to facilitate efficient experience utilization. This ensures that only high-quality experiences are retained and that the total number of experiences remains within a manageable range, preventing it from growing indefinitely with the increasing number of issues.

#### Prompt for Reflection LM: Summarize experience from the new trajectory

As a software engineering expert, you will be given an issue and summarize experiences from the resolving trajectories. Experiences usually refer to the reason for the reproduction failure and some insights. Pay special attention to this step when the model determines that the reproduction task has failed, but do not follow its format.

Repository: <repo>  
Issue: <issue>  
Trajectory: <trajectory>

When summarizing your experiences, please carefully compare the issue with the Golden Test Patch. You MUST NOT allow the newly added patch to pass the tests of your reproduction code but fail the Golden Test Patch tests, so please read carefully and summarize which parts are missing from your reproduction code.

Your output should follow the format below. Please only output the list, do not output any other text.

For all repositories:

1. ...
2. ...

For <repo>:

1. ...
2. ...

If you believe this experience is relevant to all repositories and not just limited to <repo>, please write them after 'For all repositories:'. If this experience is only applicable to the <repo> repository, write it after 'For <repo>:'. Note that the content in the two parts should NOT have any repetitions.

#### 4.4 Hierarchical Experience Pool

Based on the findings in Section 3, some error patterns are consistent across repositories, while others are repository-specific. To better maintain and utilize the extracted experiences, we design a hierarchical experience pool structure: a general pool at the top level captures common experiences shared across all repositories, followed by a range of repository-specific pools that maintain experiences unique to individual repositories. The common experiences are visible and collaboratively maintained by all issues, whereas repository-specific experiences are only visible and maintained by issues related to the specific repository. This architecture enables the model to distinguish between common and repo-specific experiences more effectively, ensures that different types of experiences

are managed and utilized effectively, thereby enhancing the overall performance and adaptability of the model.

### Prompt for Reflection LM: updating experience using multiple actions

You are an advanced reasoning agent capable of modifying your existing experiences (represented as rules) by adding, editing, removing, or merging rules based on new rules provided. Your task is to update the 'Existing rules' according to the 'New rules' and ensure the final output includes a maximum of 4 operations.

New rules:  
<new\_rules>

Existing rules:  
<existing\_rules>

You may perform the following operations:

AGREE: Choose this option if the new rules are present in the existing rules and you think they are very important.

REMOVE: Select this if the new rule contradicts existing rules or if there's redundancy among the existing rules.

ADD: Use this to introduce new rules that substantially differ from existing ones and are applicable to relevant tasks.

EDIT: Opt for this if an existing rule lacks clarity or generality. Revise it for improvement or to address past issues.

MERGE: Use this to consolidate two similar existing rules into a single, cohesive rule.

Each operation must closely follow the specified format:

<OPERATION> <RULE NUMBER>: <RULE>

The format for each operation is as follows:

AGREE <EXISTING RULE NUMBER>: <EXISTING RULE>

REMOVE <EXISTING RULE NUMBER>: <EXISTING RULE>

EDIT <EXISTING RULE NUMBER>: <NEW MODIFIED RULE>

ADD <NEW RULE NUMBER>: <NEW RULE>

MERGE <EXISTING RULE NUMBER1> <EXISTING RULE NUMBER2>: <NEW RULE>

Please follow the output format:

For all repositories:

1. ADD or EDIT or REMOVE or AGREE or MERGE ...
2. ...

For <repo>:

1. ADD or EDIT or REMOVE or AGREE or MERGE ...
2. ...

Please ensure:

1. There are no repetitions between the "For all repositories" and "For <repo>" sections.
2. If the length of EXISTING RULES is greater than 20, you must use remove or merge at least once.

Below are the operations you do to the above list of EXISTING RULES:

## 4.5 The Reproduction Process

Upon a new issue, the Actor LM selects both common and repo-specific experiences from the experience pool. For each repository, at most 10 experiences can be selected. The selected experiences are ranked by their importance. Each experience has an initial importance of 2 and is increased by 1 if it is agreed upon. The selected experiences are incorporated into the prompt for resolving the current issue.

The Reflection LM summarizes experiences from the resolving history (i.e., trajectories of issue-resolving dialogues), determines the manipulation action, and updates the experience pool.

The reproduction process continues until the model concludes that it has successfully replicated the issue. To prevent the model from prematurely declaring the task complete, we incorporate

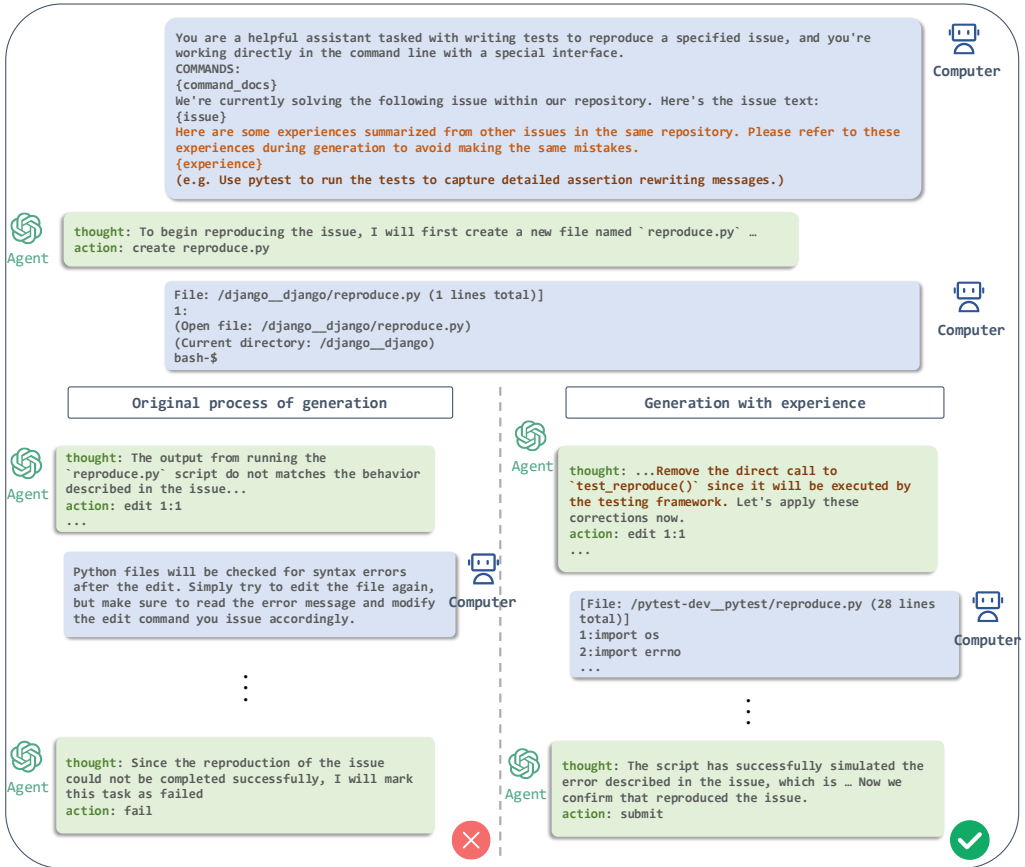


Fig. 4. Illustration of the issue reproduction process

additional verification steps. These steps involve ensuring that the model does not exhibit any of the five intrinsic errors identified in Section 3. Only when all preset conditions are satisfied is the issue considered to have been finally resolved.

An illustration example of the issue reproduction process is provided in Figure 4.

## 5 Evaluation

In this section, we present a comprehensive evaluation to assess the effectiveness of our method. We address three primary research questions (RQs).

- **RQ1: How effective is EvoCoder in issue code reproduction?** We conducted experiments to explore whether the continuous learning pipeline helps LLMs in generating reproduction code. We further conduct an ablation study of each component, in particular, whether EvoCoder can better acquire general and repository-specific knowledge, as well as learning from previous debugging processes to fix errors.

- **RQ2: Does the reproduced issue code empower issue resolving?** The reproduced code can be integrated into the issue resolving pipeline to facilitate bug localization (e.g., Spectrum-Based Fault Localization (SBFL) [1, 2]). It can also be utilized as a test script to verify the correctness of a generated function, thus facilitating code generation. This RQ aims to explore the efficacy of our method when integrated into the entire issue resolving pipeline, with a focus on the improvements in bug localization and the generation of debugging environments.
- **RQ3: What types of errors can EvoCoder help resolve?** To understand why EvoCoder is effective in resolving issues, we conduct an in-depth analysis of how it addresses the errors made by previous methods (as discussed in Section 3). Specifically, we identify and summarize the nuanced ways in which EvoCoder transforms these errors into less severe forms, thereby enabling more efficient final resolutions.

## 5.1 RQ1: Effectiveness in Issue Code Reproduction

*5.1.1 Dataset.* We utilize the widely used SWE-bench benchmark [16], which is designed to test the capability to address practical software engineering challenges. For a fair comparison, we focus on a refined subset called SWE-bench Lite [6]. This subset comprises 300 instances from SWE-bench that have been sampled to be more self-contained and covers 11 out of the original 12 repositories, ensuring a comparable range and distribution of projects as the full dataset. Our goal is to have the model attempt to reproduce and solve the specific issues described in each dataset entry.

*5.1.2 Metrics.* We did not use automatic metrics because there is a lack of ground-truth answers for reproducing code in the selected dataset. In addition, other metrics, such as the failure of reproduced code when run on the original code, do not fully equate to successful reproduction. Therefore, we ask LLM and human to score the quality of the reproduced code, adhering to five criteria: 1) the reproduction precisely aligns with the issue description; 2) the code contains no syntax or logical errors; 3) the replication process must NOT involve any form of mocking, simulation, or re-implementation of core logic that substitutes real interactions when such interactions are necessary to reproduce the issue. 4) The reproduction code should correctly interact with the necessary systems or components to produce an authentic replication of the issue. 5) the execution result of the reproduction code captures and demonstrates the key aspect of the issue as described.

1. LLM-as-a-judge scores. We leverage GPT-4 [3, 32] to judge the success of reproduced code. The input provided to the model includes the content of the issue, the reproduced code, and the execution result of the code. The LLM is asked to analyze whether the reproduced code meets the five criteria. Based on the results, the model concludes whether the code successfully reproduces the issue. The prompt used for the scoring is as follows:

### Prompt for LLM Judgement

As a software engineering expert, you are tasked with assessing the effectiveness of reproduction code in replicating an original issue as described. The input will be received in the following format:

Issue: <issue\_description>

Reproduction Code: <code>

Execution Result of Reproduction Code: <exec\_result>

Your task is to review the reproduction code and output, then answer whether the code successfully reproduced the issue, strictly avoiding any form of simulation or re-implementation of core logic that substitutes real interactions. You are required to use the following format:

Analyse from all aspects:

1. Alignment with Issue Description: Ensure the reproduction code precisely aligns with the issue described, targeting the core problem and interacting with the actual components mentioned. The code must invoke the actual methods/classes from the framework or library being discussed. But do not care about the version.
2. Code Problem Check: Ensure the code has no syntax or logical errors.
3. Avoidance of Mocking: The replication process must NOT involve any form of mocking, simulation, or

re-implementation of core logic that substitutes real interactions when such interactions are necessary to reproduce the issue. For instance, directly emulating component behavior without using the actual implementation described in the issue is not acceptable.

4. Correct Interaction: The reproduction code should correctly interact with the necessary systems or components to produce an authentic replication of the issue.

5. Demonstration of Error Cases: Confirm that the execution result of the reproduction code captures and demonstrates the key aspect of the issue as described. Pay special attention to differentiating between expected behavior as described in the issue and the actual behavior observed. The reproduction code should include tests for both valid and invalid inputs if applicable, to demonstrate any issues with error messages or unexpected behavior clearly.

Answer: [Success or Fail]  
 Error Type: [1-5] (indicate only the first encountered issue, list the corresponding number from Analysis Points, and include this only if the answer is "Fail")  
 Please conduct your analysis based on the aspects outlined above and provide a detailed answer.

2. Human scores. We ask human developers to judge the success of reproduction according to the same criteria in Section 3. In our experiment results, we observed high consistency between human and model assessments, indicating that both methods provide reliable and credible evaluation outcomes.

5.1.3 *Baselines.* We compared our method with three state-of-the-art approaches for issue code reproduction:

**SWE-agent** [44]: a system that facilitates LM agents to autonomously use computers to solve software engineering tasks, including issue resolving. SWE-agent provides an Agent-Computer Interface that includes actions such as opening and editing files and executing commands, allowing the model to interact almost freely with the computer. In the paper, the system is employed to address the full issue resolving process, with reproduction being one of its stages.

**CodeR** [8]: an issue resolving approach built upon SWE-agent. The method adopts a multi-agent architecture, designating the reproducer as a specialized agent. It also introduces format constraints for the generated reproduction code.

**LIBRO** [17]: an initial exploration into using large language models to accomplish the task of issue code reproduction. The primary method involves employing few-shot examples to prompt the LLM to generate more effective reproduction code. To adapt this method to SWE-bench, we employed a multi-round conversational approach similar to CodeR, and included two successful reproduction code snippets from the same library in the initial prompt to guide the model.

5.1.4 *Results.* As shown in Table 1, EvoCoder exhibits a significant improvement in accuracy, indicating that the continuous learning approach indeed enhances the model’s performance on similar tasks. Although few-shot prompting can be somewhat effective, its impact is limited. One reason is that identifying the most relevant examples can be challenging and may not be solely based on character or semantic similarity. Another reason is that, without explicit guidance, the model may struggle to reflect and generalize useful insights from the provided examples for subsequent reproduction tasks.

5.1.5 *Ablation Study.* We also conducted an ablation study on our method. Specifically, we tested the removal of the action-based update mechanism, the exclusion of repository-specific experiences, and the elimination of general experiences. We found that all components are crucial for the final results. The action-based mode ensures that the model’s experience is updated in a timely and flexible manner, preventing it from becoming overly bloated or rigid. Removing repo-specific experiences renders the learning process too generalized, hindering the acquisition of repository-specific knowledge. Conversely, removing general experiences isolates the experience pool within each repository, preventing the sharing of common knowledge, which can lead to a lack of experience when encountering a new repository with no existing experience.

Table 1. Comparison of reproduction accuracy by various methods

| Method                          | Accuracy (%)                 |                              |
|---------------------------------|------------------------------|------------------------------|
|                                 | LLM-judged                   | Human-judged                 |
| SWE-agent                       | 26.67                        | 27.00                        |
| CodeR                           | 34.00                        | 33.33                        |
| LIBRO                           | 45.00                        | 45.00                        |
| <b>EvoCoder (Ours)</b>          | <b>54.00</b> $\uparrow 9.00$ | <b>53.33</b> $\uparrow 8.33$ |
| - w/o Action                    | 45.66 $\downarrow 8.34$      | 46.00 $\downarrow 7.33$      |
| - w/o Repo-Specific Experiences | 46.67 $\downarrow 7.33$      | 46.33 $\downarrow 7.00$      |
| - w/o General Experiences       | 49.33 $\downarrow 4.67$      | 49.00 $\downarrow 4.33$      |

## 5.2 RQ2: Effect on Issue Resolving

We investigate the impact of our issue reproduction method in the entire issue resolving pipeline. In other words, how EvoCoder enhances the overall process of identifying, diagnosing, and fixing defects, thereby improving the quality and maintainability of software.

*5.2.1 Setup.* We take bug fix as an instance of issue resolving, where the initially generated code fails to pass the test cases. To fix the bug, the model generates a patch, applies it to the reproduced code, and verifies whether the bug has been fixed. If the patch fails, the model iterates through debugging and patch regeneration, continuing this cycle until the code is corrected or a maximum of three attempts (as set in our experiments) are reached.

We apply EvoCoder to two state-of-the-art debugging methods, AutoCodeRover [50] and Agentless [40], and compare the performance before and after applying EvoCoder. We measure the performance using the number of resolved issues and the rate of issue fixes. For the Agentless approach, we use the reproduction code to filter generated patches. We merge patches based on their occurrence frequency and select only those that produce different outputs before and after applying the patch, prioritizing the one with the highest frequency.

*5.2.2 Result.* As shown in Table 2, incorporating the reproduction code by EvoCoder has a significant impact on the model’s accuracy. For example, by applying EvoCoder to AutoCodeRover, the number of resolved issues increases from 15 to 18, an approximation of 20% improvement. The results are consistent across two basic debuggers and metrics, suggesting the effect of EvoCoder in issue resolving.

Table 2. Comparison of Reproduction Ratios

|               | # of Resolved Issues | Fix Rate (%)          |
|---------------|----------------------|-----------------------|
| AutoCodeRover | 15                   | 22.06                 |
| + w/ EvoCoder | 18 $\uparrow 3$      | 26.47 $\uparrow 4.41$ |
| Agentless     | 16                   | 23.53                 |
| + w/ EvoCoder | 18 $\uparrow 2$      | 26.47 $\uparrow 2.54$ |

## 5.3 RQ3: Error Type Transitions

To understand how EvoCoder facilitates issue resolving, we revisit the seven reproduction errors made by CodeR as discussed in Section 3. We perform a more in-depth analysis to examine how EvoCoder addresses or mitigates each of these errors.

Figure 5 shows the transition matrix of error types from CodeR and our method, along with their respective distributions. Our method successfully addresses a significant portion of issues previously encountered by CodeR. For instance, it resolves 40.9% of environment setting errors identified in CodeR.

For issues not directly resolved by our current approach, the encountered errors often shift from simpler issues, such as wrong invocation or over-mocking, to more complex challenges involving deeper logical reproduction. For example, 33.3% of the errors caused by incorrect reproduction target have been mitigated to environment setting errors. This suggests that EvoCoder effectively addresses simpler issues, primarily leaving only the more complicated ones unresolved. We also observe cases where other errors transition into incorrect reproduction targets, particularly when the issue descriptions are ambiguous. In these instances, the model may generate a response broadly relevant to the repository instead of specifically tailored to the issue, reflecting a partial understanding of the problem. For such complex cases, we anticipate that further improvements to the base model, or exploring more detailed task decomposition strategies, could yield more accurate resolutions.

Moreover, the standardized process we implemented within the Actor framework has proven especially effective in reducing issues related to environment misconfiguration and misunderstandings of instructions.

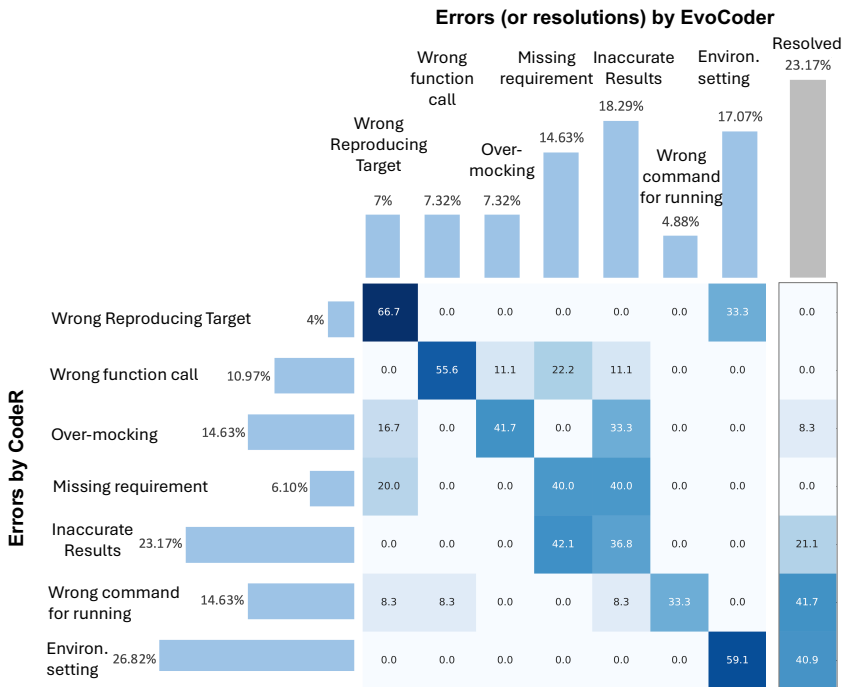


Fig. 5. The transition matrix of error types from CodeR to our method. The number in each cell denotes the percentage of errors in CodeR that has transitioned to EvoCoder

## 6 Case Study

To further verify the effectiveness of EvoCoder in real-world scenarios, we analyze two cases from SWE-bench to demonstrate the experiences we have gathered and how these experiences assist in addressing subsequent issues. The results are shown in Figures 6.

The first example pertains to an issue with the 'pytest' library. The model merely included an assert statement within the 'test\_reproduce()' function without performing any additional operations. The code is run by simply running the 'python' command when no prior experience exists. After incorporating the experience gained from previous tasks, the model did not explicitly call the defined functions within the reproduction script; instead, it allowed the 'pytest' command to automatically scan and check the tests. This approach is specific to 'pytest' and differs from the practices used with other libraries. Consequently, without specific experience with this library, the model would not be able to learn or adopt this method.

The second reproduction case involves a specific time formatting requirement within an error message. Our experiences specified that the model's output should align with the details provided in the issue, yet the model's output failed to capture the error message both before and after integrating these experiences. This situation highlights two key points: First, despite the incorporation of prior experience, the model's inherent limitations may still prevent it from successfully reproducing the issue, underscoring the need for further enhancement of the model's natural language comprehension capabilities. Second, such challenges commonly occur when the issue description is limited to textual explanations without concrete elements, such as reproducible code, expected results, or detailed specifications. For large models to effectively assist with coding issues, it is crucial to provide them with clear, comprehensive input.

## 7 Limitations and Future Work

Despite achieving significant results, our work has certain limitations, pointing to potential research directions that warrant further exploration.

First, the current research derives experiences directly from dialogue history, resulting in relatively generalized guidelines. Future work could benefit from a more granular approach—such as analyzing specific actions taken during the repair process—to generate more targeted and actionable recommendations.

Second, existing code generation [15, 27, 34, 42, 51] technologies lack comprehensive handling of boundary conditions in the way that unit tests do. As a result, even reproduced code may still exhibit issues. A promising direction for enhancement could involve integrating code generation with automated unit test generation. This combined approach could leverage the strengths of both methods while addressing their individual limitations, leading to more robust testing outcomes.

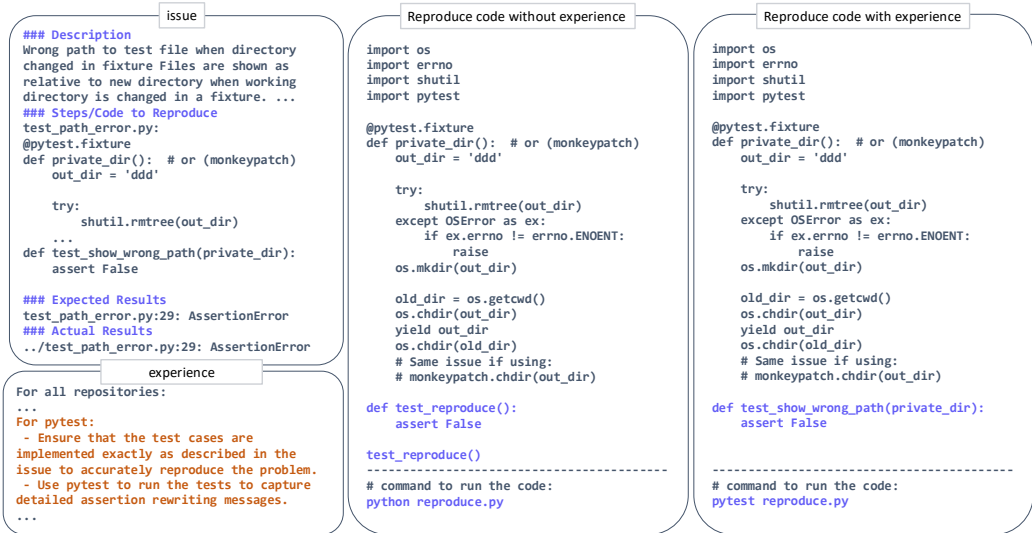
While this study primarily addresses issue code reproduction, our method holds significant potential for broader application. In future research, our approach could be extended to a range of coding scenarios, such as code translation [13, 21, 33], code editing and refactoring [4, 7, 19, 36, 47, 48]. By collaboratively building and maintaining a comprehensive repository of best practices and insights for these tasks, we can improve development efficiency and promote knowledge sharing to better support the software development community.

## 8 Related Work

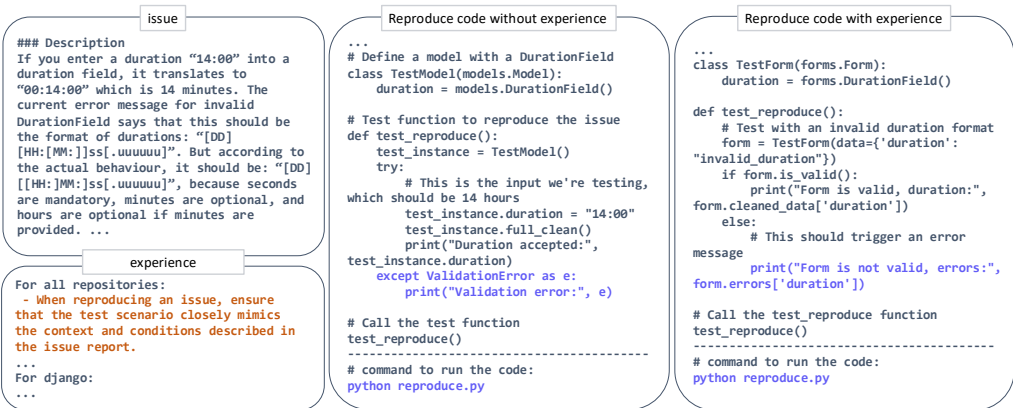
### 8.1 Repository-level Issue Resolving

With the introduction of Devin [11], the world's first AI programmer, researchers have begun to pursue the goal of enabling AI not only to assist in programming but also to independently complete software development and repair tasks [18, 22]. Towards this objective, SWE-bench [16]





(a) The reproduction code for issue 'pytest-dev\_\_pytest-7220', utilizing repository-specific experience for pytest and reproducing the issue successfully



(b) The reproduction code for issue 'django\_\_django-11049', utilizing general experience but failing to reproduce the issue.

Fig. 6. Two cases of issue reproduction by our method

serves as the first benchmark of this task. The benchmark compiles 2,294 problem-pull request pairs from 12 popular Python repositories. Unlike existing coding benchmarks (e.g. HumanEval [9], MBPP [5], CodeContests [20]), which typically focus on self-contained problems solvable within a few lines of code, SWE-bench more closely reflects the challenges encountered in real-world software development.

Based on this benchmark, numerous works have emerged to address repository-level issue resolving. One category of work involves designing a series of actions for the model, allowing it to

freely explore how to perform repairs. For example, SWE-agent [44] provides an Agent-Computer Interface (ACI) that enables the model to interact freely with the computer and complete tasks using its own capabilities. Building upon this, CodeR [8] adopts a multi-agent architecture, simulating a team to collaboratively complete tasks.

Another category of work structures the entire process into three stages: defect localization, code generation, and debugging and filtering. AutoCodeRover [50] was the first to propose this workflow, while RepoUnderstander [28] utilized Monte Carlo methods to assist in defect localization, and MarsCode Agent [24] employed code knowledge graphs and language server protocols to enhance the model's understanding of the code repository structure.

Our work primarily focuses on generating reproducible code for issue replication, with the primary aim of providing executable code to facilitate the model's debugging process and is also designed to be readily integrable with other methodologies.

## 8.2 Issue Code Reproduction

Issue code reproduction is attracting increasing attention from researchers. LIBRO [17] was the first to propose using large models to accomplish the task of reproducing problematic code, prompting these models with a few examples to generate more effective reproduction code.

With the development of various agent strategies [11, 24, 28, 38, 40, 44, 50], it has been found that adopting a multi-round dialogue approach can more effectively activate the capabilities of models, thereby achieving better results. Many solutions addressing the full-chain issues of agents in SWE-Bench [16] have already incorporated this element [8, 44]. However, they have only added some fixed tips in a relatively simplistic manner, without implementing more refined designs.

Furthermore, some research attempts to directly generate corresponding unit tests from issue reports [30]. The test functions produced by this method can integrate better into the overall content of the repository. However, the generation of unit tests may be a more challenging task, as it requires simultaneous consideration of code location and generation while ensuring that the addition of new tests does not disrupt the existing structure. Observations indicate that the accuracy of unit test generation is currently lower than that of the aforementioned code reproduction methods.

In this study, we further improve the success rate of code reproduction based on the Agent-Computer Interface provided by SWE-agent [44] and the mechanism of multi-round dialogues. Meanwhile, we view the automatic generation of unit tests as a more ambitious and challenging long-term goal, anticipating a transition towards this advanced form of reproduction when model capabilities significantly improve in the future.

## 9 Conclusion

In this paper, we propose a novel continuous learning framework for issue code reproduction. Our method maintains a hierarchical pool for both common and repo-specific experiences. By continuously resolving new issues, the model accumulates experiences specific to each repository and progressively updates its common knowledge base. Experimental results demonstrate that our method achieves a 20% improvement in issue code reproduction compared to the state-of-the-art methods. Additionally, our approach significantly enhances performance in the repository-level issue resolving tasks.

## 10 Data Availability Statement

We release our code to encourage further exploration in this direction. The artifact that supports the results discussed in this paper is available at <https://anonymous.4open.science/r/EvoCoder-6433/>

## References

- [1] Rui Abreu, Peter Zoetewij, Rob Golsteijn, and Arjan JC Van Gemund. 2009. A practical evaluation of spectrum-based fault localization. *Journal of Systems and Software* 82, 11 (2009), 1780–1792.
- [2] Rui Abreu, Peter Zoetewij, 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.
- [3] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- [4] Eman Abdullah AlOmar, Anushkrishna Venkatakrishnan, Mohamed Wiem Mkaouer, Christian Newman, and Ali Ouni. 2024. How to refactor this code? An exploratory study on developer-ChatGPT refactoring conversations. In *Proceedings of the 21st International Conference on Mining Software Repositories*. 202–206.
- [5] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732* (2021).
- [6] Carlos E. Jimenez, John Yang, Jiayi Geng. 2024. *SWE-bench Lite: A Canonical Subset for Efficient Evaluation of Language Models as Software Engineers*. <https://www.swebench.com/lite.html>
- [7] Saikat Chakraborty and Baishakhi Ray. 2021. On multi-modal learning of editing source code. In *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 443–455.
- [8] Dong Chen, Shaoxin Lin, Muhan Zeng, Daoguang Zan, Jian-Gang Wang, Anton Cheshkov, Jun Sun, Hao Yu, Guoliang Dong, Artem Aliev, et al. 2024. CodeR: Issue Resolving with Multi-Agent and Task Graphs. *arXiv preprint arXiv:2406.01304* (2024).
- [9] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating Large Language Models Trained on Code. (2021). *arXiv:2107.03374 [cs.LG]*
- [10] Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2023. Teaching large language models to self-debug. *arXiv preprint arXiv:2304.05128* (2023).
- [11] Cognition. 2023. *Introducing Devin*. <https://www.cognition.ai/introducing-devin>
- [12] Yangruibo Ding, Marcus J Min, Gail Kaiser, and Baishakhi Ray. 2024. Cycle: Learning to self-refine the code generation. *Proceedings of the ACM on Programming Languages* 8, OOPSLA1 (2024), 392–418.
- [13] Linyuan Gong, Mostafa Elhoushi, and Alvin Cheung. 2024. AST-T5: Structure-Aware Pretraining for Code Generation and Understanding. *arXiv preprint arXiv:2401.03003* (2024).
- [14] Xiangbing Huang, Yingwei Ma, Haifang Zhou, Zhijie Jiang, Yuanliang Zhang, Teng Wang, and Shanshan Li. 2023. Towards Better Multilingual Code Search through Cross-Lingual Contrastive Learning. In *Proceedings of the 14th Asia-Pacific Symposium on Internetware*. 22–32.
- [15] Zhijie Jiang, Haixu Xiong, Yingwei Ma, Yao Zhang, Yan Ding, Yun Xiong, and Shanshan Li. 2023. Automatic Code Annotation Generation Based on Heterogeneous Graph Structure. In *2023 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 497–508.
- [16] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770* (2023).
- [17] Sungmin Kang, Juyeon Yoon, and Shin Yoo. 2023. Large language models are few-shot testers: Exploring llm-based general bug reproduction. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2312–2323.
- [18] Jiaying Li, Yan Lei, Shanshan Li, Haifang Zhou, Yue Yu, Zhouyang Jia, Yingwei Ma, and Teng Wang. 2023. A two-stage framework for ambiguous classification in software engineering. In *2023 IEEE 34th International Symposium on Software Reliability Engineering (ISSRE)*. IEEE, 275–286.
- [19] Jia Li, Ge Li, Zhuo Li, Zhi Jin, Xing Hu, Kechi Zhang, and Zhiyi Fu. 2023. Codeeditor: Learning to edit source code with pre-trained models. *ACM Transactions on Software Engineering and Methodology* 32, 6 (2023), 1–22.
- [20] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz,

- Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. 2022. Competition-level code generation with AlphaCode. *Science* 378, 6624 (2022), 1092–1097. <https://doi.org/10.1126/science.abq1158> arXiv:<https://www.science.org/doi/pdf/10.1126/science.abq1158>
- [21] Bingchang Liu, Chaoyu Chen, Zi Gong, Cong Liao, Huan Wang, Zhichao Lei, Ming Liang, Dajun Chen, Min Shen, Hailian Zhou, et al. 2024. Mftcoder: Boosting code llms with multitask fine-tuning. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 5430–5441.
- [22] Junwei Liu, Kaixin Wang, Yixuan Chen, Xin Peng, Zhenpeng Chen, Lingming Zhang, and Yiling Lou. 2024. Large Language Model-Based Agents for Software Engineering: A Survey. *arXiv preprint arXiv:2409.02977* (2024).
- [23] Wenzhuo Liu, Fei Zhu, and Cheng-Lin Liu. 2024. Branch-Tuning: Balancing Stability and Plasticity for Continual Self-Supervised Learning. *arXiv preprint arXiv:2403.18266* (2024).
- [24] Yizhou Liu, Pengfei Gao, Xinchun Wang, Chao Peng, and Zhao Zhang. 2024. MarsCode Agent: AI-native Automated Bug Fixing. *arXiv preprint arXiv:2409.00899* (2024).
- [25] Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evol-instruct. *arXiv preprint arXiv:2306.08568* (2023).
- [26] Yingwei Ma, Rongyu Cao, Yongchang Cao, Yue Zhang, Jue Chen, Yibo Liu, Yuchen Liu, Binhua Li, Fei Huang, and Yongbin Li. 2024. Lingma SWE-GPT: An Open Development-Process-Centric Language Model for Automated Software Improvement. *arXiv preprint arXiv:2411.00622* (2024).
- [27] Yingwei Ma, Yue Liu, Yue Yu, Yuanliang Zhang, Yu Jiang, Changjian Wang, and Shanshan Li. 2023. At Which Training Stage Does Code Data Help LLMs Reasoning? *arXiv preprint arXiv:2309.16298* (2023).
- [28] Yingwei Ma, Qingping Yang, Rongyu Cao, Binhua Li, Fei Huang, and Yongbin Li. 2024. How to Understand Whole Software Repository? *arXiv preprint arXiv:2406.01422* (2024).
- [29] Yingwei Ma, Yue Yu, Shanshan Li, Zhouyang Jia, Jun Ma, Rulin Xu, Wei Dong, and Xiangke Liao. 2023. Mulcs: Towards a unified deep representation for multilingual code search. In *2023 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 120–131.
- [30] Niels Müндler, Mark Niklas Müller, Jingxuan He, and Martin Vechev. 2024. Code Agents are State of the Art Software Testers. *arXiv preprint arXiv:2406.12952* (2024).
- [31] Ansong Ni, Miltiadis Allamanis, Arman Cohan, Yinlin Deng, Kensen Shi, Charles Sutton, and Pengcheng Yin. 2024. NExT: Teaching Large Language Models to Reason about Code Execution. *arXiv preprint arXiv:2404.14662* (2024).
- [32] OpenAI. 2024. *Introducing GPT-4o*. <https://openai.com/index/hello-gpt-4o/>
- [33] Rangeet Pan, Ali Reza Ibrahimzada, Rahul Krishna, Divya Sankar, Lambert Pouguem Wassi, Michele Merler, Boris Sobolev, Raju Pavuluri, Saurabh Sinha, and Reyhaneh Jabbarvand. 2023. Understanding the effectiveness of large language models in code translation. *arXiv preprint arXiv:2308.03109* (2023).
- [34] Zhenyu Pan, Rongyu Cao, Yongchang Cao, Yingwei Ma, Binhua Li, Fei Huang, Han Liu, and Yongbin Li. 2024. Codev-Bench: How Do LLMs Understand Developer-Centric Code Completion? *arXiv preprint arXiv:2410.01353* (2024).
- [35] Yuling Shi, Songsong Wang, Chengcheng Wan, and Xiaodong Gu. 2024. From Code to Correctness: Closing the Last Mile of Code Generation with Hierarchical Debugging. *arXiv preprint arXiv:2410.01215* (2024).
- [36] Atsushi Shirafuji, Yusuke Oda, Jun Suzuki, Makoto Morishita, and Yutaka Watanobe. 2023. Refactoring programs using large language models with few-shot examples. In *2023 30th Asia-Pacific Software Engineering Conference (APSEC)*. IEEE, 151–160.
- [37] Bo Wang, Tianxiang Sun, Hang Yan, Siyin Wang, Qingyuan Cheng, and Xipeng Qiu. 2024. In-Memory Learning: A Declarative Learning Framework for Large Language Models. *arXiv preprint arXiv:2403.02757* (2024).
- [38] Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. 2024. Executable Code Actions Elicit Better LLM Agents. arXiv:2402.01030 [cs.CL]
- [39] Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. 2024. Magicoder: Empowering code generation with oss-instruct. In *Forty-first International Conference on Machine Learning*.
- [40] Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2024. Agentless: Demystifying llm-based software engineering agents. *arXiv preprint arXiv:2407.01489* (2024).
- [41] Xin Xia, David Lo, Xinyu Wang, and Bo Zhou. 2013. Accurate developer recommendation for bug resolution. In *2013 20th Working Conference on Reverse Engineering (WCRE)*. IEEE, 72–81.
- [42] Ruiyang Xu, Jialun Cao, Yaojie Lu, Hongyu Lin, Xianpei Han, Ben He, Shing-Chi Cheung, and Le Sun. 2024. CRUXEval-X: A Benchmark for Multilingual Code Reasoning, Understanding and Execution. *arXiv preprint arXiv:2408.13001* (2024).
- [43] 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).
- [44] John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. Swe-agent: Agent-computer interfaces enable automated software engineering. *arXiv preprint arXiv:2405.15793* (2024).

- [45] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629* (2022).
- [46] Zhaojian Yu, Xin Zhang, Ning Shang, Yangyu Huang, Can Xu, Yishujie Zhao, Wenxiang Hu, and Qiufeng Yin. 2023. Wavocoder: Widespread and versatile enhanced instruction tuning with refined data generation. *arXiv preprint arXiv:2312.14187* (2023).
- [47] Mengxiao Zhang, Yongqiang Tian, Zhenyang Xu, Yiwen Dong, Shin Hwei Tan, and Chengnian Sun. 2023. Lampr: Boosting the Effectiveness of Language-Generic Program Reduction via Large Language Models. *arXiv preprint arXiv:2312.13064* (2023).
- [48] Mengxiao Zhang, Yongqiang Tian, Zhenyang Xu, Yiwen Dong, Shin Hwei Tan, and Chengnian Sun. 2024. LPR: Large Language Models-Aided Program Reduction. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*. 261–273.
- [49] Tao Zhang, He Jiang, Xiapu Luo, and Alvin TS Chan. 2016. A literature review of research in bug resolution: Tasks, challenges and future directions. *Comput. J.* 59, 5 (2016), 741–773.
- [50] 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.
- [51] Qiming Zhu, Jialun Cao, Yaojie Lu, Hongyu Lin, Xianpei Han, Le Sun, and Shing-Chi Cheung. 2024. DOMAINEVAL: An Auto-Constructed Benchmark for Multi-Domain Code Generation. *arXiv preprint arXiv:2408.13204* (2024).
- [52] Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, Huazuo Gao, Shirong Ma, et al. 2024. DeepSeek-Coder-V2: Breaking the Barrier of Closed-Source Models in Code Intelligence. *arXiv preprint arXiv:2406.11931* (2024).
- [53] Novak Zuber, Gary E Wilson, Mamoru Ishii, Wolfgang Wulff, BE Boyack, AE Dukler, P Griffith, JM Healzer, RE Henry, JR Lehner, et al. 1998. An integrated structure and scaling methodology for severe accident technical issue resolution: development of methodology. *Nuclear Engineering and Design* 186, 1-2 (1998), 1–21.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009