# InfCode: Adversarial Iterative Refinement of Tests and Patches for Reliable Software Issue Resolution

KEFAN LI, Beihang University, China and Beijing Tokfinity Technology Co., Ltd., China

MENGFEI WANG, HENGZHI ZHANG, and ZHICHAO LI, Beijing Tokfinity Technology Co., Ltd., China

YUAN YUAN, MU LI, XIANG GAO, HAILONG SUN, CHUNMING HU, and WEIFENG LV, Beihang University, China

Large language models have advanced software engineering automation, yet resolving real-world software issues remains difficult because it requires repository-level reasoning, accurate diagnostics, and strong verification signals. Existing agent-based and pipeline-based methods often rely on insufficient tests, which can lead to patches that satisfy verification but fail to fix the underlying defect. We present **InfCode**, an adversarial multi-agent framework for automated repository-level issue resolution. InfCode iteratively refines both tests and patches through adversarial interaction between a Test Patch Generator and a Code Patch Generator, while a Selector agent identifies the most reliable fix. The framework runs inside a containerized environment that supports realistic repository inspection, modification, and validation. Experiments on SWE-bench Lite and SWE-bench Verified using models such as DeepSeek-V3 and Claude 4.5 Sonnet show that InfCode consistently outperforms strong baselines. It achieves **79.4%** performance on SWE-bench Verified, establishing a new state-of-the-art. We have released InfCode as an open-source project at https://github.com/Tokfinity/InfCode.

## 1 Introduction

Large language models (LLMs) have recently demonstrated strong capabilities across many software engineering tasks. Their ability to interpret natural language specifications, generate executable code, and interact with external tools has motivated increasing interest in automating software issue resolution. This task requires diagnosing the behavioral discrepancy described in an issue report and synthesizing a correct code modification that integrates into a complex repository. Software issue resolution refers to identifying the root cause of a reported bug and producing an appropriate code modification that restores the intended behavior, ensure the reliable of software systems [12, 15, 16]. Compared with isolated code tasks, resolving repository-level issues requires a more comprehensive understanding of multi-file dependencies, project-specific invariants, and execution behaviors.

LLMs have achieved impressive results at the function level in code generation [5, 8, 20], code repair [31, 37], and test generation [17, 30], yet repository-level issue resolution remains substantially

Authors' Contact Information: Kefan Li, Beihang University, Beijing, China and Beijing Tokfinity Technology Co., Ltd., Beijing, China; Mengfei Wang; Hengzhi Zhang; Zhichao Li, Beijing Tokfinity Technology Co., Ltd., Beijing, China; Yuan Yuan; Mu Li; Xiang Gao; Hailong Sun; Chunming Hu; Weifeng Lv, Beihang University, Beijing, China.

more challenging. To mitigate these difficulties, two main lines of work have emerged. Agent-based approaches such as SWE-agent [44], OpenHands [40], Moatless-Tools [46], AutoCodeRover [45], SpecRover [36], and Trae-Agent [13] provide LLMs with tool interfaces for repository exploration, including file editing, shell execution, and code searching. These systems benefit from flexible and adaptive reasoning, as agents can iteratively gather information and validate intermediate hypotheses. Pipeline-based methods, such as Agentless [42], decompose the task into structured subtasks with targeted prompting, which improves stability and reduces error accumulation during reasoning.

Recent systems often attempt to generate tests that reproduce the reported issue, moving beyond reliance on existing repository tests. While this represents an important advancement, the generated tests are frequently not sufficiently strong, not accurately aligned with the issue semantics, or unable to capture subtle behavioral constraints. As a result, agents may optimize code patches against weak or imperfect tests and thus fail to fully resolve the underlying problem. Moreover, most prior work treats test generation and code generation as loosely connected steps. Without a mechanism that explicitly coordinates the two, the verification signal remains limited and may permit patches that only satisfy partial or over-simplified test conditions.

To address these limitations, we propose **InfCode**, an adversarial multi-agent framework for automated repository-level issue resolution. The framework introduces two specialized agents that engage in iterative adversarial refinement. A Test Patch Generator constructs and strengthens test cases based on the issue description, aiming to expose the incorrect behavior more effectively. A Code Patch Generator responds to these strengthened tests by producing improved code modifications. This adversarial interaction drives both agents toward more rigorous testing and more robust patches. A Selector agent subsequently evaluates all candidate patches and identifies the most reliable one. All agents operate inside a containerized environment that ensures reproducible execution and provides a suite of tools for repository inspection, modification, and validation.

We evaluate InfCode on SWE-bench Lite [16] and DeepSeek-V3 [21] and further test it with stronger models such as Claude 4.5 Sonnet [3] on SWE-bench Verified [33]. The results show consistent improvements over strong agent-based and pipeline-based baselines and establish a new state-of-the-art (SOTA) (79.4%) on SWE-bench Verified. Ablation studies confirm the effectiveness of adversarial iteration and the importance of the final selection stage. This demonstrates that adversarially iterating between test generation and code patching provides a principled and effective strategy for enhancing the reliability and generality of automated software issue resolution. In summary, we propose the following contributions:

- We propose an adversarial multi-agent framework that iteratively refines both code and test patches through generation and selection, leading to higher patch quality and robustness.
- We develop a repository-aware tool suite that enhances agents' understanding and manipulation of real-world codebases within a containerized environment.
- We conduct comprehensive experiments showing that our approach outperforms existing automated patch generation methods in both correctness and generalization.

## 2 Method

### 2.1 Framework Overview

Our framework adopts a multi-agent architecture that automates software patch generation and validation through adversarial collaboration between agents. As illustrated in Figure 1, the overall process is divided into two stages: Stage 1: Patch Generation and Stage 2: Patch Selection. The system begins with the automatic construction of a controlled environment using a containerized image builder, followed by the orchestration of specialized tools, such as bash utilities, code editors,

searchers, submitters, and executors, that support the end-to-end patch generation and evaluation workflow.
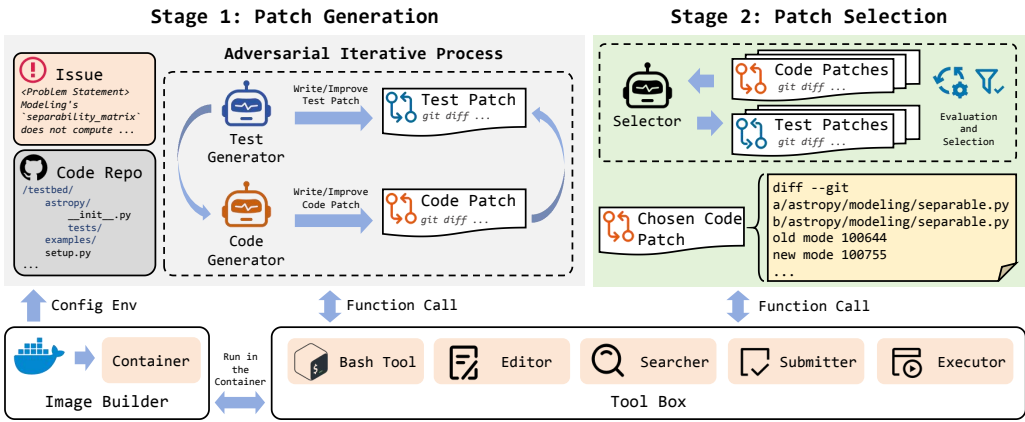


Fig. 1. Overview of InfCode for automated code patch generation and selection.

At the core of the framework lies a dual-agent adversarial generation mechanism, where the Code Patch Generator Agent and the Test Patch Generator Agent interact iteratively to produce and refine code patches. The primary motivation for this design is that existing test cases often fail to expose the problematic behavior described in a reported issue. Therefore, we introduce an adversarial loop in which the Test Generator continuously strengthens its test suite, while the Code Generator incrementally refines the code to satisfy increasingly stringent requirements. This interplay ensures that the final code patch is not only functional but also robust against a comprehensive set of test conditions.

By combining containerized reproducibility, multi-agent cooperation, and adversarial iteration, the framework enables the automatic synthesis of high-quality code patches that effectively resolve reported issues while maintaining software reliability.

## 2.2 Stage 1: Patch Generation

Stage 1 focuses on producing candidate patches through an **adversarial iterative process** between the Test Generator and the Code Generator. Given an issue and its associated code repository, the Test Generator first creates test patches that attempt to reproduce the faulty behavior. The Code Generator then produces code patches aimed at passing these newly generated tests. Both agents operate inside the controlled container environment, invoking necessary operations through the tool suite.

*Adversarial Iterative Refinement.* The core innovation of this stage lies in the adversarial iterative refinement mechanism between code and test generation. When the Code Generator produces a patch that successfully passes the current test suite, the Test Generator re-analyzes the issue to identify weaknesses, missing edge cases, or insufficient coverage. It then introduces additional or stronger test cases that challenge the existing implementation. In response, the Code Generator must refine its code to satisfy these enhanced tests. This iterative competition and cooperation between the two agents leads to a dynamic equilibrium where both the tests and code progressively improve: tests become increasingly comprehensive and discriminative, while code becomes more robust and generalizable. Ultimately, this adversarial generation strategy encourages the emergence

of high-quality patches that address the issue thoroughly and resist regression under unseen conditions. To prevent unbounded iteration between the Test Generator and the Code Generator, we impose a hard cap on the number of iterations, once this limit is reached, the generation process terminates immediately. In addition, if the Test Generator strengthens the test suite and the code still passes all tests, the refinement loop terminates at once.

## 2.3    Stage 2: Patch Selection

After a series of candidate patches are produced in Stage 1, Stage 2 focuses on identifying the optimal code patch through systematic evaluation and selection. The Selector agent collects all generated code patches along with their associated test patches and evaluates them based on predefined performance metrics such as functional correctness, test coverage, execution success, and compatibility with repository constraints.

The Selector may execute the candidate patches within the containerized environment to verify their behavior empirically. It then compares their performance, filtering out overfitted or unstable solutions. Ultimately, the Selector selects the most reliable and generalizable code patch—the one that not only resolves the reported issue but also passes the strengthened test suite produced during adversarial iteration. This chosen patch represents the final output of the framework and can be directly integrated into the target repository.

## 2.4    Tool Suite and Execution Environment

To ensure reproducibility and consistent execution, all operations in our framework are performed within a Docker container instantiated for each issue. The container provides an isolated environment where a suite of tools collaboratively supports patch generation, editing, and evaluation. The following tools are integrated within the containerized workflow:

- **Bash Tool:** Executes shell commands within the container, enabling task automation such as running tests, managing dependencies, and invoking build scripts.
- **Editor:** Supports file creation, content insertion and replacement, and retrieval of specific line ranges, allowing fine-grained source code modification.
- **Searcher:** Performs efficient iterative directory searches using `ripgrep`, with support for regular expressions to locate relevant code or test fragments.
- **Submitter:** Runs `git diff` to extract and submit patch contents, maintaining version control consistency and recording patch provenance.
- **Executor:** Acts as the interface between tools and the Docker container, forwarding execution requests and managing input/output interactions.

This tool suite collectively enables automated, reproducible, and isolated code patch generation and evaluation within the adversarial multi-agent framework.

## 3   Experimental Setup

### 3.1   Research Questions

We propose the following research questions (RQs):

- **RQ1:** How well does the proposed agent framework perform in generating correct and robust patches compared with existing methods?
- **RQ2:** How does each component of the framework contribute to its overall effectiveness?
- **RQ3:** What is the maximum performance achieved by InfCode when applied to SOTA LLMs?

## 3.2 Benchmarks

We conducted the experiments of RQ1 using SWE-bench Lite [16], a lightweight subset of the SWE-bench dataset, which is particularly suitable for lightweight evaluations. Many baselines have been evaluated on this subset. To optimize resource usage, we employed the DeepSeek-V3 model for the experiments. DeepSeek-V3 [21] possesses strong code generation and tool invocation capabilities, making it an ideal choice as the backbone LLM for agent systems. Additionally, to assess the absolute performance of our method, we evaluated it on the SWE-bench Verified [33] subset and the highly capable model Claude 4.5 Sonnet [3]. The SWE-bench Verified subset is a curated version of the SWE-bench dataset, where both solutions and tests have undergone rigorous manual validation. The tests in this subset are sufficiently rigorous, enabling a more stringent evaluation of patches. Our method achieved SOTA performance on the SWE-bench Verified leaderboard.

## 3.3 Metrics

We employed the evaluation method provided by the official SWE-bench. After generating the patches, we used the `git diff` command to extract them. It is important to note that, in order to avoid interference from modifications to test files, we excluded all files starting with `test` from the evaluation. Subsequently, we conducted the evaluation using the SWE-bench command-line interface. We reported the solved rate, the number of problems solved, and the average cost per problem (in USD).

## 3.4 Baselines

On SWE-bench Lite, we compared various baselines that performed well on this dataset. The solved rate and cost data were obtained from their respective papers/reports or the SWE-bench Lite leaderboard. To ensure a fair comparison, we primarily selected experimental data from models such as DeepSeek-V3 [21] and GPT-4o [32]. On SWE-bench Verified, we compared our method with the top five performing methods, with the data sourced from the SWE-bench Verified leaderboard.

## 4 Results

## 4.1 RQ1: Overall Effectiveness Comparison

*Performance on SWE-bench Lite.* Table 1 presents the performance of InfCode compared to other baseline methods on the SWE-bench Lite dataset. We observe that, among the 300 problems in the SWE-Bench Lite dataset, InfCode with Deepseek-V3 as the backbone model successfully solves 118 problems. This performance is the best among baselines utilizing similar models. Furthermore, although our method incurs higher costs than certain baselines, such as the KGCompass, it outperforms them. Additionally, thanks to the relatively low API cost of the Deepseek-V3 model, our approach significantly reduces costs while improving performance, compared to methods using GPT-4.

*Unique issues fixed.* Figure 2 illustrates the problem coverage and differences between InfCode and the baseline methods. Among the problems solved, 36 are common to all methods. In addition, among all the methods, InfCode and KGCompass solved the most unique problems. Furthermore, InfCode solved 11 more problems than KGCompass. This highlights the superior performance of the InfCode method.

Table 1. Comparison of InfCode and other baseline methods on the SWE-bench Lite dataset.

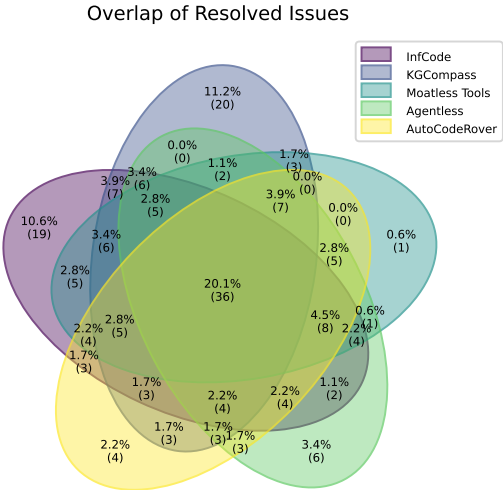| Method | LLM | SWE-bench Lite | |
|---|---|---|---|
| | | Resolved (%) | Avg. Cost ($) |
| SWE-agent [44] | GPT-4o | 18.33% (55) | 2.53 |
| OpenHands [40] | GPT-4o | 22.00% (66) | 1.72 |
| Agentless [42] | GPT-4o | 32.00% (96) | 0.70 |
| SpecRover [36] | Sonnet-3.5+GPT-4o | 31.00% (93) | 0.65 |
| AutoCodeRover [45] | GPT-4 | 19.00% (57) | 0.43 |
| | GPT-4o | 30.67% (92) | - |
| SWE-Search [4] | GPT-4o | 31.00% (93) | - |
| Moatless Tools [46] | DeepSeek-V3 | 30.67% (92) | - |
| | GPT-4o | 24.67% (74) | - |
| KGCompass [24] | DeepSeek-V3 | 36.67% (110) | 0.20 |
| **InfCode** | **DeepSeek-V3** | **40.33% (121)** | 0.26 |



Fig. 2. Illustrate the overlap between issues resolved by InfCode and the top baseline methods.

Answer RQ1: On the SWE-bench Lite dataset, using the DeepSeek-V3 model and models of similar capabilities, InfCode solved the most problems and also addressed the highest number of unique problems, highlighting the superior performance of the InfCode method.

## 4.2 RQ2: Ablation Study of Framework Components

To investigate the contribution of each module in InfCode, we conducted ablation experiments. Specifically, we removed the Adversarial Iteration and Selection modules, resulting in the configurations "w/o Adversarial" and "w/o Selection", and performed experiments on the SWE-bench Lite dataset using the DeepSeek-V3 model. The results are presented in Table 2.

Table 2. Ablation study of InfCode on the SWE-bench Lite dataset using DeepSeek-V3.

| Method | Resolved (%) |
|---|---|
| **InfCode** | **40.33% (121)** |
| w/o Adversarial | 36.33% (109) |
| w/o Selection | 32.33% (97) |

The results demonstrate that both modules contribute to the performance of InfCode, with the removal of either module leading to a decrease in performance. Notably, the Selection module has a greater impact on the overall performance than the Adversarial Iteration module. This highlights the effectiveness of both the adversarial iteration and patch selection strategies.

> Answer RQ2: Both the Adversarial Iteration and Selection modules contribute to the performance of InfCode, with the Selection module making a larger contribution.

## 4.3 RQ3: Performance on SOTA LLMs

Although the performance on SWE-bench Lite is promising, it may suffer from a potential issue of insufficient test intensity, which could lead to an inability to filter out potentially incorrect patch repairs. Additionally, the current DeepSeek-V3 [21] model is not the strongest model, and performance on the most advanced models would be more persuasive. Therefore, we also conducted experiments on the SWE-bench Verified [33] dataset using the strongest model, Claude 4.5 Sonnet [3]. The results are presented in Table 3.

Table 3. Comparison of InfCode (based on the Claude 4.5 Sonnet) with the top 5 methods on the SWE-bench Verified leaderboard as of November 14, 2025.

| Method | Resolved (%) | Rank |
|---|---|---|
| 🏆 **InfCode** | **79.4% (397)** | **1** |
| TRAE + Doubao-Seed-Code | 78.80% (394) | 2 |
| Atlassian Rovo Dev (2025-09-02) | 76.80% (384) | 3 |
| EPAM AI/Run Developer Agent v20250719 + Claude 4 Sonnet | 76.80% (384) | 3 |
| ACoder | 76.40% (382) | 5 |
| Warp | 75.60% (378) | 6 |

In Table 3, we report the performance of InfCode on Claude 4.5 Sonnet, along with the performance of the top five methods, excluding InfCode, on the current SWE-bench Verified Leaderboard. The results show that InfCode achieved an impressive pass rate, securing the top position. This further confirms the effectiveness of InfCode.

Answer RQ3: On the SWE-bench Verified dataset, InfCode with Claude 4.5 Sonnet as the backbone model achieves the top performance, demonstrating the effectiveness of InfCode.

## 5 Discussion

In this section, we analyze the failure rate of tool invocations by LLMs. The system provides four tools that can be called directly by the model: Bash Tool, Editor, Searcher, and Submitter. We collected statistics on the average number of calls and the failure rates of these tools when InfCode solves a problem. The distribution of call counts for each tool is shown in the left panel of Figure 3, and the failure rates for each tool are reported in the right panel of Figure 3.
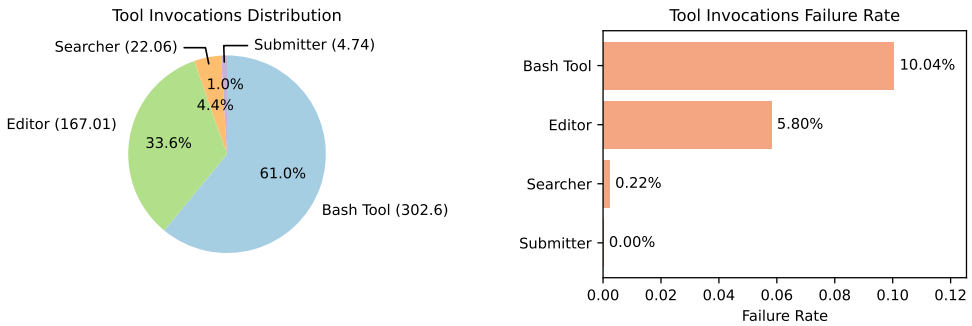


Fig. 3. Distribution of the average number of tool invocations per problem and the corresponding failure rates.

From the distribution of tool invocations, InfCode calls the Bash Tool most frequently, followed by the Editor. This indicates that most attempts focus on modifying files and executing code. InfCode performs extensive file inspection, code editing, and debugging, and these operations contribute substantially to successful problem solving. The next most common is the Searcher, which primarily supports initial problem localization, once the relevant files are identified, the Editor is used to inspect their contents. The Submitter is invoked least often and is used only when the model judges that the problem has been solved. From the failure rate plot, the Bash Tool exhibits the highest failure rate. Inspection of execution logs shows that the main causes include attempts to run complex scripts with `python -c "xxx"` and attempts to execute nonexistent commands, largely due to unfamiliarity with the environment or available commands. The Editor has the next highest failure rate, mostly because its string replacement function requires an exact `old_str`, which can be challenging in some contexts. Overall, the failure rates across tools remain low, which demonstrates the robustness of InfCode.

## 6 Limitations

Although InfCode attains strong performance, several limitations and avenues for improvement remain. First, prior works [1, 2, 28] have shown that LLM-based agents may produce incorrect test patches when reproducing issues, which can misdirect subsequent code repair. Our log analysis reveals a similar pattern: during adversarial iterations, the Test Generator sometimes creates specialized tests that deviate from the issue description in order to induce failures, which in turn misleads the Code Generator. This suggests that strategies are still needed to improve the accuracy and faithfulness of the Test Generator. Second, the Bash Tool and the Editor continue to exhibit

some invocation errors. The implementations of these tools require further refinement, and the usage guidelines should be clarified to reduce failure rates as much as possible.

## 7 Threads to Validity

*Internal Validity.* One potential internal threat arises from the possibility that the generated patches might inadvertently modify the corresponding test files, which could bias the evaluation results. Such modifications may artificially increase the apparent success rate of generated patches by weakening or bypassing test cases. To mitigate this risk, we explicitly exclude all files with names beginning with "test" when computing the patch differences using the `git diff` command. This ensures that only genuine code modifications are considered in the evaluation, preserving the integrity and fairness of the results.

*External Validity.* A potential external threat concerns the generalizability of our approach to programming languages other than those evaluated in our experiments. While our current study focuses on a specific language environment, the proposed framework itself is designed to be language-agnostic. It operates primarily through regular-expression-based search, command-line execution, and repository-level operations driven by bash commands. Consequently, with minor adaptations to language-specific syntax and build tools, the framework can be readily extended to other programming languages and ecosystems.

*Construct Validity.* Construct validity relates to whether our evaluation metrics and experimental setup accurately reflect the true effectiveness of automated patch generation. A possible concern is that passing existing or generated test cases may not fully represent real-world correctness or semantic equivalence to developer fixes. To address this, we evaluate patches not only by test success but also by robustness under adversarially strengthened test suites and by cross-validation against multiple independent test generations. This design ensures that our measurements align with the intended construct of patch quality and that the evaluation outcomes genuinely reflect functional correctness and generalizability.

## 8 Related Work

### 8.1 Repository-Level Issue Resolution

Repository-level software issue resolution has emerged as a central challenge in leveraging LLMs for real-world development workflows. The introduction of SWE-bench [16] provided an executable benchmark grounded in authentic GitHub issues with verifiable test-based evaluation, highlighting a substantial gap between the impressive surface-level generation capabilities of LLMs and the robustness required for practical software maintenance. Subsequent studies have shown that even reported "solved" issues are frequently incorrect or incomplete [41], underscoring the need for both reliable evaluation and improved reasoning mechanisms.

To bridge this gap, multiple efforts explore agent-based approaches for repository navigation, reasoning, and patch synthesis. Frameworks such as SWE-agent [44], OpenHands [40], and MarsCode [22] offer general-purpose agent execution platforms, while systems like AutoCodeRover [45] and CodeR [7] introduce multi-agent or role-specialized collaboration guided by task graphs or iterative retrieval. LingmaAgent [25] extends repository exploration to include historical commits, pull requests, and dependency relationships, and Lingma SWE-GPT [23] integrates development-process-centric knowledge to improve coherence with real engineering workflows. Meanwhile, AGENTLESS [42] questions the necessity of complex agent orchestration, demonstrating that reducing step-wise error propagation can improve repair reliability.

Another line of work focuses on enhancing repository understanding and structured reasoning. RepoGraph [34] and repository-aware knowledge graph approaches [43] inject structural code relationships, while SpecRover [36] targets intent extraction to improve semantic grounding. Search-enhanced methods such as SWE-Search [4] employ Monte Carlo Tree Search for guided repair, and competitive strategies like SWE-Debate [18] use multi-agent argumentation to refine solutions. Experience-driven systems, including SWE-Exp [9] and EXPEREPAIR [27], incorporate memory mechanisms to transfer prior solution strategies across repositories.

Recent findings further demonstrate that compute scaling during inference, rather than simply model size, is crucial for successful repository-level repair. Thinking Longer, Not Larger [24] and Trae Agent [13] show that increased test-time reasoning depth and reflection significantly improve performance. Complementary work such as BugPilot [38] focuses on generating complex bug scenarios to better train and stress-test these systems.

## 8.2 Iterative Code–Test Feedback

Iterative refinement driven by execution feedback has emerged as a crucial paradigm for improving the correctness and reliability of LLM-generated code. Early work on model self-debugging [10] demonstrated that LLMs can leverage compilation or runtime error messages to analyze failures and propose targeted revisions, establishing the principle that feedback grounded in actual program behavior is essential for meaningful correction. This idea was further generalized in Self-Refine [26], which introduced a task-agnostic iterative feedback loop in which model outputs are repeatedly critiqued and refined. Building upon these concepts, Reflexion [37] incorporates verbal self-reflection to enable agents to accumulate reasoning strategies across iterations, while RLEF [14] employs reinforcement learning from execution performance signals to directly shape the model's repair policy.

A complementary direction enhances the feedback loop by involving both code and tests. LLM-LOOP [35] and CoCoEvo [19] jointly generate tests and candidate patches, fostering co-evolution of validation and repair, whereas LEVER [29] emphasizes verification-in-the-loop, using constraint checks and invariant validation to filter incorrect solutions. More recent work revisits self-debugging through self-generated tests [11], demonstrating that adaptive test synthesis substantially improves fault localization. Meanwhile, RepairAgent [6] operationalizes iterative feedback within an autonomous agent framework, integrating code navigation, environment interaction, and repeated refinement to solve real repair tasks.

Beyond system design, theoretical analyses show that iterative code repair inherently involves an exploration–exploitation trade-off [39], where excessive refinement may lead to local minima, while insufficient refinement may fail to converge. Together, these works converge on the insight that effective code repair relies not only on initial generation quality, but also on structured, feedback-driven iterative improvement tightly aligned with executable program semantics.

## 9 Conclusion

This paper introduced InfCode, an adversarial multi-agent framework designed to improve the reliability and robustness of automated repository-level issue resolution. The framework integrates iterative collaboration between a Test Patch Generator and a Code Patch Generator, enabling tests to be progressively strengthened while patches are refined to satisfy increasingly rigorous requirements. This adversarial co-evolution produces verification signals that more accurately capture the true semantics of reported issues. In addition, the Selector agent evaluates candidate patches within a controlled containerized environment, ensuring consistent execution, reducing instability, and identifying the most reliable solution.

Extensive experiments on SWE-bench Lite and SWE-bench Verified demonstrate that InfCode achieves substantial improvements over existing agent-based and pipeline-based approaches. With Claude 4.5 Sonnet as the backbone model, InfCode attains a solve rate of 79.4% on SWE-bench Verified, which represents the current state of the art. Ablation studies further highlight the contributions of adversarial iteration and patch selection, confirming that both components are essential for maximizing system performance. Overall, InfCode shows that coupling adversarial test refinement with iterative patch generation provides an effective strategy for advancing automated software repair. Future research may explore applying this framework to additional programming languages and development ecosystems, enhancing the intelligence of test generation, and integrating learned heuristics to guide interaction between agents more efficiently.

# References

[1] Toufique Ahmed, Jatin Ganhotra, Rangeet Pan, Avraham Shinnar, Saurabh Sinha, and Martin Hirzel. 2025. Otter: Generating Tests from Issues to Validate SWE Patches. *arXiv preprint arXiv:2502.05368* (2025).

[2] Toufique Ahmed, Martin Hirzel, Rangeet Pan, Avraham Shinnar, and Saurabh Sinha. 2024. TDD-Bench Verified: Can LLMs Generate Tests for Issues Before They Get Resolved? *arXiv preprint arXiv:2412.02883* (2024).

[3] Anthropic. 2025. *Claude Sonnet 4.5.* Accessed: 2025-11-13.

[4] Antonis Antoniades, Albert Örwall, Kexun Zhang, Yuxi Xie, Anirudh Goyal, and William Wang. 2024. Swe-search: Enhancing software agents with monte carlo tree search and iterative refinement. *arXiv preprint arXiv:2410.20285* (2024).

[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] Islem Bouzenia, Premkumar Devanbu, and Michael Pradel. 2024. Repairagent: An autonomous, llm-based agent for program repair. *arXiv preprint arXiv:2403.17134* (2024).

[7] 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).

[8] Mark Chen. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).

[9] Silin Chen, Shaoxin Lin, Xiaodong Gu, Yuling Shi, Heng Lian, Longfei Yun, Dong Chen, Weiguo Sun, Lin Cao, and Qianxiang Wang. 2025. Swe-exp: Experience-driven software issue resolution. *arXiv preprint arXiv:2507.23361* (2025).

[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] Xiancai Chen, Zhengwei Tao, Kechi Zhang, Changzhi Zhou, Wanli Gu, Yuanpeng He, Mengdi Zhang, Xunliang Cai, Haiyan Zhao, and Zhi Jin. 2025. Revisit self-debugging with self-generated tests for code generation. *arXiv preprint arXiv:2501.12793* (2025).

[12] Yoav Freund and Robert E Schapire. 1997. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences* 55, 1 (1997), 119–139.

[13] Pengfei Gao, Zhao Tian, Xiangxin Meng, Xinchen Wang, Ruida Hu, Yuanan Xiao, Yizhou Liu, Zhao Zhang, Junjie Chen, Cuiyun Gao, et al. 2025. Trae agent: An llm-based agent for software engineering with test-time scaling. *arXiv preprint arXiv:2507.23370* (2025).

[14] Jonas Gehring, Kunhao Zheng, Jade Copet, Vegard Mella, Quentin Carbonneaux, Taco Cohen, and Gabriel Synnaeve. 2024. Rlef: Grounding code llms in execution feedback with reinforcement learning. *arXiv preprint arXiv:2410.02089* (2024).

[15] Lianghong Guo, Wei Tao, Runhan Jiang, Yanlin Wang, Jiachi Chen, Xilin Liu, Yuchi Ma, Mingzhi Mao, Hongyu Zhang, and Zibin Zheng. 2025. Omnigirl: A multilingual and multimodal benchmark for github issue resolution. *Proceedings of the ACM on Software Engineering* 2, ISSTA (2025), 24–46.

[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] Han Li, Yuling Shi, Shaoxin Lin, Xiaodong Gu, Heng Lian, Xin Wang, Yantao Jia, Tao Huang, and Qianxiang Wang. 2025. Swe-debate: Competitive multi-agent debate for software issue resolution. *arXiv preprint arXiv:2507.23348* (2025).

[19] Kefan Li, Yuan Yuan, Hongyue Yu, Tingyu Guo, and Shijie Cao. 2025. CoCoEvo: Co-Evolution of Programs and Test Cases to Enhance Code Generation. *IEEE Transactions on Evolutionary Computation* (2025).

[20] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. 2022. Competition-level code generation with alphacode. *Science* 378, 6624 (2022), 1092–1097.

[21] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437* (2024).

[22] 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).

[23] 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).

[24] Yingwei Ma, Yongbin Li, Yihong Dong, Xue Jiang, Rongyu Cao, Jue Chen, Fei Huang, and Binhua Li. 2025. Thinking longer, not larger: Enhancing software engineering agents via scaling test-time compute. *arXiv preprint arXiv:2503.23803* (2025).

[25] Yingwei Ma, Qingping Yang, Rongyu Cao, Binhua Li, Fei Huang, and Yongbin Li. 2025. Alibaba lingmaagent: Improving automated issue resolution via comprehensive repository exploration. In *Proceedings of the 33rd ACM International Conference on the Foundations of Software Engineering*. 238–249.

[26] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2023. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems* 36 (2023), 46534–46594.

[27] Fangwen Mu, Junjie Wang, Lin Shi, Song Wang, Shoubin Li, and Qing Wang. 2025. EXPEREPAIR: Dual-Memory Enhanced LLM-based Repository-Level Program Repair. *arXiv preprint arXiv:2506.10484* (2025).

[28] Niels Mündler, Mark Müller, Jingxuan He, and Martin Vechev. 2024. SWT-bench: Testing and validating real-world bug-fixes with code agents. *Advances in Neural Information Processing Systems* 37 (2024), 81857–81887.

[29] Ansong Ni, Srini Iyer, Dragomir Radev, Veselin Stoyanov, Wen-tau Yih, Sida Wang, and Xi Victoria Lin. 2023. Lever: Learning to verify language-to-code generation with execution. In *International Conference on Machine Learning*. PMLR, 26106–26128.

[30] Pengyu Nie, Rahul Banerjee, Junyi Jessy Li, Raymond J Mooney, and Milos Gligoric. 2023. Learning deep semantics for test completion. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2111–2123.

[31] Theo X Olausson, Jeevana Priya Inala, Chenglong Wang, Jianfeng Gao, and Armando Solar-Lezama. 2023. Is self-repair a silver bullet for code generation? *arXiv preprint arXiv:2306.09896* (2023).

[32] OpenAI. 2024. GPT-4o. https://openai.com. Large language model.

[33] OpenAI. 2024. Introducing SWE-bench Verified. https://openai.com/index/introducingâĂŚsweâĂŚbenchâĂŚverified/. Accessed: YYYY-MM-DD.

[34] 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).

[35] Ravin Ravi, Dylan Bradshaw, Stefano Ruberto, Gunel Jahangirova, and Valerio Terragni. 2025. LLMLOOP: Improving LLM-Generated Code and Tests through Automated Iterative Feedback Loops. *ICSME. IEEE* (2025).

[36] Haifeng Ruan, Yuntong Zhang, and Abhik Roychoudhury. 2024. Specrover: Code intent extraction via llms. *arXiv preprint arXiv:2408.02232* (2024).

[37] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems* 36 (2023), 8634–8652.

[38] Atharv Sonwane, Isadora White, Hyunji Lee, Matheus Pereira, Lucas Caccia, Minseon Kim, Zhengyan Shi, Chinmay Singh, Alessandro Sordoni, Marc-Alexandre Côté, et al. 2025. BugPilot: Complex Bug Generation for Efficient Learning of SWE Skills. *arXiv preprint arXiv:2510.19898* (2025).

[39] Hao Tang, Keya Hu, Jin Zhou, Si Cheng Zhong, Wei-Long Zheng, Xujie Si, and Kevin Ellis. 2024. Code repair with llms gives an exploration-exploitation tradeoff. *Advances in Neural Information Processing Systems* 37 (2024), 117954–117996.

[40] Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. 2024. Openhands: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741* (2024).

[41] You Wang, Michael Pradel, and Zhongxin Liu. 2025. Are" Solved Issues" in SWE-bench Really Solved Correctly? An Empirical Study. *arXiv preprint arXiv:2503.15223* (2025).

[42] Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2024. Agentless: Demystifying llm-based software engineering agents. *arXiv preprint arXiv:2407.01489* (2024).

[43] Boyang Yang, Haoye Tian, Jiadong Ren, Shunfu Jin, Yang Liu, Feng Liu, and Bach Le. 2025. Enhancing repository-level software repair via repository-aware knowledge graphs. *arXiv preprint arXiv:2503.21710* (2025).

[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. *Advances in Neural Information Processing Systems* 37 (2024), 50528–50652.

[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] Albert Örwall. 2024. Moatless Tools. https://github.com/aorwall/moatless-tools. Accessed: 2024-11-13.