

# Revisit Self-Debugging with Self-Generated Tests for Code Generation

Anonymous ACL submission

## Abstract

Large language models (LLMs) have demonstrated significant advancements in code generation, yet they still face challenges when tackling tasks that extend beyond their basic capabilities. Recently, the concept of self-debugging has been proposed as a way to enhance code generation performance by leveraging execution feedback from tests. However, the availability of high-quality tests in real-world scenarios is often limited. In this context, self-debugging with self-generated tests emerges as a promising solution, though its limitations and practical potential have not been fully explored. To address this gap, we investigate the efficacy of self-debugging in code generation tasks. We propose and analyze two distinct paradigms for the self-debugging process: post-execution and in-execution self-debugging. Our findings reveal that post-execution self-debugging struggles with the test bias introduced by self-generated tests, which can lead to misleading feedback. In contrast, in-execution self-debugging enables LLMs to mitigate this bias and leverage intermediate states during program execution. By focusing on runtime information rather than relying solely on potentially flawed self-generated tests, this approach demonstrates significant promise for improving the robustness and accuracy of LLMs in code generation tasks.

## 1 Introduction

Large language models (LLMs) have made significant strides in the field of code generation, but still face challenges to perform complex programming tasks beyond their basic capabilities (Jain et al., 2024). Such tasks demand LLMs to comprehend requirements in natural language and produce programs that behave as expected. Recently, Self-Debugging has emerged as a promising method to boost the performance of LLMs in code generation (Chen et al., 2024b; Jiang et al., 2023; Zhong

et al., 2024), which allows models to iteratively refine their outputs through execution with *pre-defined oracle tests*. A critical limitation arises in real-world software development scenarios: those oracle tests are not readily available for each programming task, posing a significant challenge for the practical application of this method.

To address this challenge, recent studies have introduced *self-generated tests* into the self-debugging process (Shinn et al., 2023; Huang et al., 2023; Ridnik et al., 2024). Beyond leveraging intrinsic capabilities (Huang et al., 2024; Madaan et al., 2023; Chen et al., 2024b), execution feedback from *self-generated tests* could also serve as additional signals to help LLMs identify defects in their code generation. This approach helps reduce reliance on external feedback from humans or stronger models. However, it remains unclear whether self-debugging using self-generated tests is a reliable and scalable solution to enhance code generation. Reflexion (Shinn et al., 2023) leverages feedback from generated tests to debug the code but triggers the repair using hidden oracle tests. AlphaCodium (Ridnik et al., 2024) allows LLMs to generate tests with supervision of public oracle tests. The improvements observed using oracle tests do not accurately demonstrate the true self-debugging capabilities of LLMs. This highlights the need for a more comprehensive and transparent evaluation to better understand the inherent debugging potential with execution feedback from tests.

In this work, we formally define two distinct paradigms for doing this: **post-execution and in-execution self-debugging**. The conventional post-execution method operates binary feedback from self-testing—pass or fail labels obtained through testing with self-generated tests, susceptible to test bias as models must infer defects solely from potentially flawed tests. We propose in-execution self-debugging, a novel paradigm that allows models to analyze intermediate states during program execu-

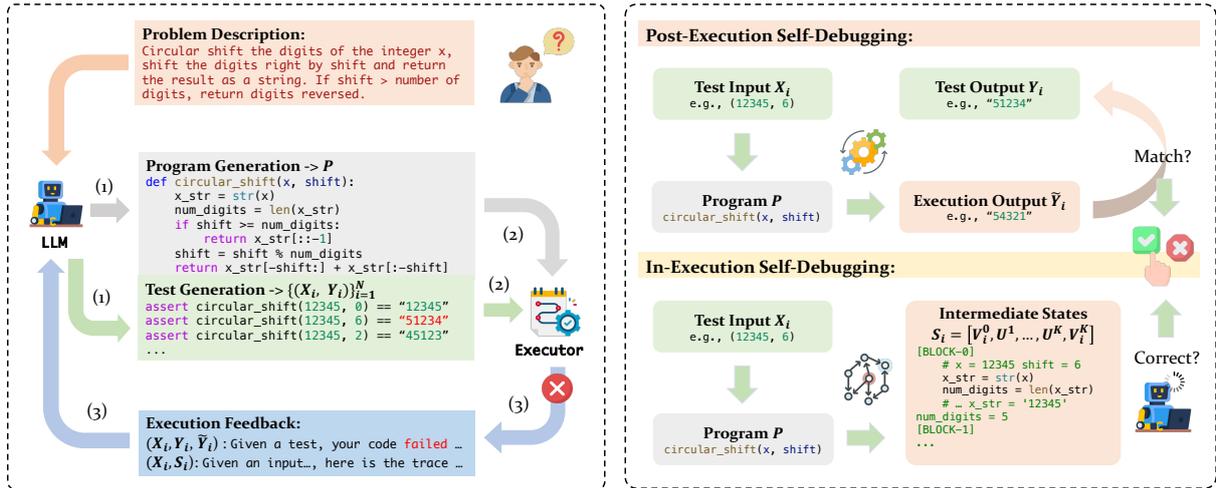


Figure 1: **Overview of self-debugging with execution feedback from self-generated tests.** (1) The model generates an initial program along with a suite of tests, based on the problem description. (2) The program is executed by an executor on the self-generated tests. (3) The feedback from execution is then utilized by the model to produce a revised version of the program.

tion without dependency on post-execution signals. This highlights a key gap: existing methods conflate in-execution traces with post-execution feedback (Zhong et al., 2024; Ni et al., 2024; Bouzenia et al., 2023), potentially failing to isolate the model from unreliable test verdicts.

We investigate the efficacy of both paradigms with self-generated tests applied to four frontier LLMs. Experimental results across both basic and competitive code generation tasks reveal key insights. We find that the post-execution method struggles with simpler tasks but shows potential for improvement on more challenging ones. The inconsistency across problem levels is attributed to the bias introduced by self-generated tests, which refers to the misalignment between self-testing labels and true labels of the programs, resulting in a critical bottleneck for the post-execution method. In-execution self-debugging transcends this limitation by solely focusing on the intermediate states during the program execution, reducing bias susceptibility. It consistently outperforms post-execution method, demonstrating its promise for advancing self-debugging capabilities in LLMs.

Our work is the first to systematically dissect the self-testing bias in self-debugging and propose a paradigm shift from post-execution label validation to in-execution state reasoning, thereby reducing dependency on unreliable test outputs. We summarize our contributions as follows:

- We formally define and evaluate two distinct paradigms: post-execution and in-execution

self-debugging, providing a comprehensive analysis of their strengths and limitations.

- We identify the critical issue of test bias in self-generated tests, demonstrating how it undermines post-execution methods.
- We propose in-execution self-debugging as a robust alternative, leveraging intermediate runtime states to reduce bias and improve debugging accuracy.

Through our study, we aim to shed light on the practicality of self-debugging with self-generated tests, contributing valuable insights into the future development of LLMs in code generation tasks.

## 2 Self-Debugging with Self-Generated Tests

We focus on evaluating the self-debugging capabilities of LLMs through execution on their self-generated tests. Figure 1 provides a comprehensive overview of this process. Given a problem with a natural language specification, the language model (denoted as  $M$ ) generates an initial program  $P$  along with a suite of tests  $\{(X_i, Y_i)\}_{i=1}^N$ , where  $X_i$  represents the input and  $Y_i$  represents the expected output for the  $i$ -th test. To enhance the model’s debugging performance beyond its intrinsic reasoning capability, we utilize execution feedback as an additional signal to help the model identify defects in its generated program according to the problem description. Specifically, we employ

an executor (denoted as E) to run the program on the test suite and collect execution information as feedback. This feedback is then used to guide the model in refining the program.

There are various implementations for utilizing execution feedback, which we categorize into two distinct paradigms: **post-execution and in-execution self-debugging**. These paradigms differ in the type of information employed during the self-debugging process. Post-execution information refers to content obtained after the program’s execution, such as execution outputs or error messages. In contrast, in-execution information refers to content observed during the program’s execution, providing finer-grained insights into its behavior. We now formally define these paradigms.

**Post-Execution Self-Debugging.** The paradigm leverages information obtained after the actual execution of the program. A widely adopted implementation involves comparing the execution output with the expected output (Olausson et al., 2024; Wang et al., 2024; Dong et al., 2023; Madaan et al., 2023; Zhang et al., 2023a; Chen et al., 2024b; Jiang et al., 2024), as illustrated in Figure 1. Given an initial program  $P$  and a set of generated tests  $\{(X_i, Y_i)\}_{i=1}^N$ , an executor E processes each input  $X_i$ , yielding the corresponding execution output:

$$\tilde{Y}_i = E(P, X_i), i \in [1, N] \quad (1)$$

The executor automatically assesses whether the execution output  $\tilde{Y}_i$  aligns with the expected output  $Y_i$  to determine the success of the  $i$ -th test. If a discrepancy or error is detected, the test is marked as failed. The system then utilizes the failed test  $(X_i, Y_i)$ , the execution output  $\tilde{Y}_i$ , and any associated error messages to repair the program. This process encourages the model M to generate a revised version of the program, denoted as:

$$\tilde{P} = M(P, X_i, Y_i, \tilde{Y}_i) \quad (2)$$

**In-Execution Self-Debugging.** Post-execution self-debugging typically overlooks the runtime information of the program, which can provide valuable insights for program repair and enhancement. To address this limitation, in-execution self-debugging leverages feedback from the intermediate states during program execution (Nye et al., 2021; Ni et al., 2024; Bouzenia et al., 2023). Formally, a program’s execution trace  $T_i$  can be divided into multiple basic units, denoted as  $T_i =$

$[U^1, U^2, \dots, U^K]$ , where  $U^k$  represents the  $k$ -th unit and  $K$  is the total number of units in the execution trace. Each unit is defined as a sequence of statements and expressions with a single entry and a single exit point. Given an input  $X_i$ ,  $i \in [1, N]$ , the executor E initializes the input as the variable set  $V_i^0 = X_i$  and passes it through the first unit  $U^1$ . The execution updates the variable set to  $V_i^1 = E(U^1, V_i^0)$ , where  $V_i^1$  denotes the set of variables after executing unit  $U^1$ . This process is repeated iteratively, with the executor traversing each subsequent unit  $U^k$  until the program execution is complete:

$$V_i^k = E(U^k, V_i^{k-1}), k \in [1, K] \quad (3)$$

The sequence of intermediate states represented as  $S_i = [V_i^0, U^1, V_i^1, \dots, U^K, V_i^K]$ , provides a detailed view of how the program behaves over time. By analyzing this runtime information, the model M is able to identify potential issues within specific units and refine the program accordingly. Notably, in-execution self-debugging excludes any post-execution signals (e.g., negative labels and corresponding test outputs), focusing solely on the dynamic behavior of the program during execution. This process results in the revised version of the program, denoted as:

$$\tilde{P} = M(P, X_i, S_i) \quad (4)$$

### 3 Experiments

In this section, we comprehensively evaluate the self-debugging capabilities of a wide range of LLMs using self-generated tests on code generation benchmarks covering diverse levels of programming problems. We carry out experiments to answer the following research questions:

**RQ1:** *Are model-generated tests reliable enough? How does self-debugging perform leveraging post-execution information from these tests?*

**RQ2:** *Is the performance of post-execution self-debugging consistent on more complex tasks? If not, what factors contribute to the discrepancy?*

**RQ3:** *How does in-execution self-debugging perform compared with the post-execution method?*

#### 3.1 Experimental Setup

Producing high-quality tests poses significant challenges for LLMs, as it necessitates a comprehensive understanding of natural language specifications and advanced code reasoning capabilities (Gu et al.,

2024; Chen et al., 2024a). Therefore, we investigate the research questions with four advanced chat models: **LLaMA-3-70B-Instruct** (Dubey et al., 2024), **Qwen2.5-Coder-7B-Instruct** (Hui et al., 2024), the API-served **GPT-4o-2024-05-13** and **Claude-3.5-Sonnet**. We select three code generation benchmarks—**HumanEval+** (Chen et al., 2021; Liu et al., 2023), **MBPP+** (Austin et al., 2021; Liu et al., 2023), and **LiveCodeBench** (Jain et al., 2024), covering both basic and competitive<sup>1</sup> programming problems to comprehensively evaluate the efficacy of self-debugging.

We employ a greedy decoding strategy (with a temperature of zero) across all generation phases of self-debugging to ensure deterministic and consistent results. We design prompts for the initial program generation to ensure that no additional information is introduced by subsequent prompts for program repair. This premise is crucial for us to concentrate on investigating the true self-debugging capabilities of LLMs (Huang et al., 2024). To generate a test suite for each problem, we prompt the model to produce  $n = 10$  diverse and extensive tests based on the corresponding natural language specification in a zero-shot manner. We discuss the effect of the varying  $n$  in Appendix D. For a detailed overview of the prompts used in our experiments, please refer to Appendix G.

### 3.2 RQ1: Post-Execution Self-Debugging Struggles with Potentially Flawed Tests

In this subsection, we evaluate the model-generated tests on basic tasks, and then examine the performance of self-debugging utilizing post-execution information from these tests, which is consistent with the conventional implementations. The correctness of a program is determined by comparing the actual output with the expected output for each test case. If the generated program manages to pass all tests, the iterative process terminates, and no further debugging is conducted.

**LLMs Cannot Yet Produce Reliable Test Suites Even for Basic Tasks.** To better understand the reliability of tests generated by the model itself, we employ program contracts and canonical solutions provided by HumanEval and MBPP to evaluate the validity of test inputs and outputs respectively. For detailed illustration for the validation, please refer

<sup>1</sup>In this work, we regard problems in HumanEval, MBPP as basic programming problems, and those in LiveCodeBench as competitive ones according to overall complexity and difficulty.

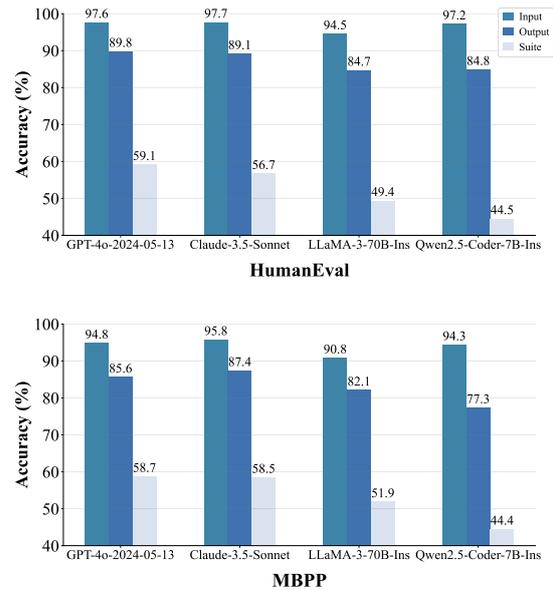


Figure 2: Accuracy of self-generated tests on HumanEval and MBPP. Test **input & output** are evaluated case-by-case; A test **suite** is deemed valid if all outputs within the suite are correct.

to Appendix C. We also calculate the overall accuracy for the test suites. Figure 2 summarizes the results. GPT-4o and Claude-3.5-Sonnet demonstrate superior capability in producing high-quality tests compared to others, yet they remain prone to generating flawed tests. While models can often generate plausible inputs, they struggle to accurately predict the corresponding outputs, which are critical for validating program behavior. Notably, there is a significant gap between the accuracy of the test suites and the accuracy when we consider the tests individually. This suggests that post-execution information from self-generated tests has a high chance of providing unreliable feedback to the self-debugging process.

**Post-Execution Self-Debugging Struggles on Basic Tasks.** To assess the impact of the potentially flawed tests, we consider two different types of feedback derived from post-execution results. The first type is the binary correct *label*, which simply indicates the correctness of the previous program. If the program is incorrect, an instruction for repair is provided to the model. The second type is the *detail* of the failure, which includes the test input, expected output, and execution output. In cases where the program raises an exception during execution, the error message is incorporated into the detail in place of the execution output. We conduct experiments on problems from the Hu-

Table 1: Results of **post-execution self-debugging** (post-exec.) with *self-generated tests*. The values highlighted in **red** are declines while those in **green** are increases relative to the initial generation (w/o *self-debug*).

Model	Method	#Iter	HumanEval+	MBPP+	LiveCodeBench			
					easy	medium	hard	overall
GPT-4o-2024-05-13	w/o <i>self-debug</i>	0	87.8	76.5	89.3	33.1	6.0	46.0
	post-exec. w/ <i>label</i>	1	87.2 <sup>-0.6</sup>	76.7 <sup>+0.2</sup>	89.9 <sup>+0.6</sup>	41.1 <sup>+8.0</sup>	6.0 <sup>+0.0</sup>	49.3 <sup>+3.3</sup>
		2	86.6 <sup>-1.2</sup>	77.5 <sup>+1.0</sup>	89.9 <sup>+0.6</sup>	40.0 <sup>+6.9</sup>	6.9 <sup>+0.9</sup>	49.1 <sup>+3.1</sup>
	post-exec. w/ <i>detail</i>	1	84.1 <sup>-3.7</sup>	76.2 <sup>-0.3</sup>	85.5 <sup>-3.8</sup>	36.0 <sup>+2.9</sup>	8.6 <sup>+2.6</sup>	46.4 <sup>+0.4</sup>
2		85.4 <sup>-2.4</sup>	76.5 <sup>+0.0</sup>	87.4 <sup>-1.9</sup>	38.3 <sup>+5.2</sup>	8.6 <sup>+2.6</sup>	48.0 <sup>+2.0</sup>	
Claude-3.5-Sonnet	w/o <i>self-debug</i>	0	89.0	77.0	93.1	48.0	16.4	55.8
	post-exec. w/ <i>label</i>	1	88.4 <sup>-0.6</sup>	77.8 <sup>+0.8</sup>	89.9 <sup>-3.2</sup>	49.1 <sup>+1.1</sup>	17.2 <sup>+0.8</sup>	55.3 <sup>-0.5</sup>
		2	86.6 <sup>-2.4</sup>	76.2 <sup>-0.8</sup>	91.2 <sup>-1.9</sup>	49.7 <sup>+1.7</sup>	16.4 <sup>+0.0</sup>	55.8 <sup>+0.0</sup>
	post-exec. w/ <i>detail</i>	1	81.1 <sup>-7.9</sup>	72.8 <sup>-4.2</sup>	89.9 <sup>-3.2</sup>	49.1 <sup>+1.1</sup>	13.8 <sup>-2.6</sup>	54.4 <sup>-1.2</sup>
2		79.3 <sup>-9.7</sup>	75.4 <sup>-1.6</sup>	85.5 <sup>-7.6</sup>	43.3 <sup>-4.7</sup>	8.6 <sup>-7.8</sup>	49.3 <sup>-6.5</sup>	
LLaMA-3-70B-Ins	w/o <i>self-debug</i>	0	73.8	71.2	72.3	10.3	2.6	30.2
	post-exec. w/ <i>label</i>	1	65.2 <sup>-8.6</sup>	68.3 <sup>-2.9</sup>	66.0 <sup>-6.3</sup>	9.1 <sup>-1.2</sup>	3.4 <sup>+0.8</sup>	27.8 <sup>-2.4</sup>
		2	69.5 <sup>-4.3</sup>	68.3 <sup>-2.9</sup>	64.8 <sup>-7.5</sup>	10.9 <sup>+0.6</sup>	2.6 <sup>+0.0</sup>	27.8 <sup>-2.4</sup>
	post-exec. w/ <i>detail</i>	1	66.5 <sup>-7.3</sup>	64.8 <sup>-6.4</sup>	56.6 <sup>-15.7</sup>	10.9 <sup>+0.6</sup>	4.3 <sup>+1.7</sup>	25.3 <sup>-4.9</sup>
2		67.1 <sup>-6.7</sup>	63.8 <sup>-7.4</sup>	63.5 <sup>-8.8</sup>	12.0 <sup>+1.7</sup>	2.6 <sup>+0.0</sup>	27.8 <sup>-2.4</sup>	
Qwen2.5-Coder-7B-Ins	w/o <i>self-debug</i>	0	81.7	70.6	74.8	23.4	8.6	35.8
	post-exec. w/ <i>label</i>	1	78.0 <sup>-3.7</sup>	69.8 <sup>-0.8</sup>	69.8 <sup>-5.0</sup>	24.0 <sup>+0.6</sup>	8.6 <sup>+0.0</sup>	34.2 <sup>-1.6</sup>
		2	79.3 <sup>-2.4</sup>	69.8 <sup>-0.8</sup>	71.7 <sup>-3.1</sup>	23.4 <sup>+0.0</sup>	8.6 <sup>+0.0</sup>	34.7 <sup>-1.1</sup>
	post-exec. w/ <i>detail</i>	1	76.2 <sup>-5.5</sup>	68.0 <sup>-2.6</sup>	69.2 <sup>-5.6</sup>	20.0 <sup>-3.4</sup>	8.6 <sup>+0.0</sup>	32.4 <sup>-3.4</sup>
2		75.6 <sup>-6.1</sup>	69.0 <sup>-1.6</sup>	66.7 <sup>-8.1</sup>	21.1 <sup>-2.3</sup>	8.6 <sup>+0.0</sup>	32.0 <sup>-3.8</sup>	

manEval and MBPP benchmarks. Table 1 presents the pass rates achieved through post-execution self-debugging, showcasing notable declines in performance on both benchmarks for LLaMA-3-70B-Instruct and Qwen2.5-Coder-7B-Instruct. For GPT-4o and Claude-3.5-Sonnet, a consistent decrease is observed on HumanEval. While the performance on MBPP may show initial improvement, the use of more detailed feedback ultimately leads to worse outcomes compared to the initial generation. The results highlight the challenges and limitations of the post-execution method with self-generated tests, even for state-of-the-art models.

**Failure Case Study** In post-execution settings, erroneous test outputs introduce ambiguity into the self-debugging process. We illustrate this challenge with an example from HumanEval using GPT-4o in Figure 7 in Appendix F. *When a test fails, the model is expected to determine whether the failure stems from defects in the program or errors in the test itself.* This uncertainty complicates the self-debugging process and necessitates a further investigation into the effects of testing on self-generated tests. While post-execution information with self-generated tests is leveraged, self-debugging remains a bottleneck, limiting im-

provements beyond initial generation.

### 3.3 RQ2: Bias from Self-Testing Leads to Inconsistency Across Levels of Tasks

To investigate the performance of self-debugging on more challenging tasks, we further conduct post-execution self-debugging experiments using problems from LiveCodeBench. The problems are classified into three distinct difficulty levels: easy, medium, and hard. We report the pass rate achieved at each level of difficulty, as well as the overall performance across the benchmark.

**Label Feedback Provides Some Relief on More Challenging Tasks.** Table 1 also summarizes the results of post-execution self-debugging on LiveCodeBench. We find that self-debugging using label feedback generally leads to improvements on medium and hard-level problems. This is notably in contrast to their performance on basic and easy problems. All models except GPT-4o experience a decrease in overall performance, primarily due to significant declines in easy problems. Despite incorporating more post-execution results, the overall performance with detailed feedback remains inferior to that achieved with label feedback.

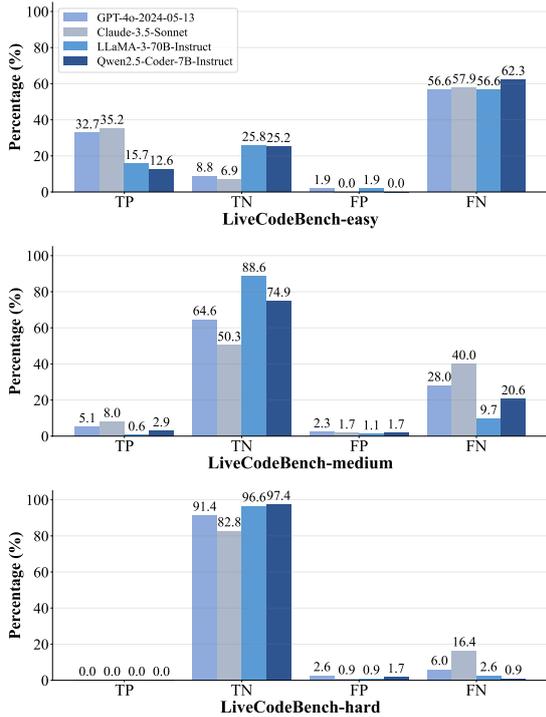


Figure 3: Test bias in labels when evaluating the programs with self-generated tests on different difficulty levels of problems in LiveCodeBench.

**Self-Testing Bias Analysis** To investigate the reasons behind the inconsistent results across different levels of problems, we delve into the impact on testing programs with self-generated tests. We acknowledge that even advanced LLMs are likely to generate inaccurate tests (Gu et al., 2024). As a result, a program that is actually correct might fail some of the generated tests, resulting in a false negative label (FN). Conversely, a flawed program might pass all the test cases, leading to a false positive label (FP), which could prevent necessary updates and prematurely present a buggy program as correct. *The misalignment between self-testing labels and true labels highlights the bias introduced by self-generated tests in program evaluation.*

We present an analysis of test bias in labels at the first iteration of self-debugging, as illustrated in Figure 3. Given the implementation of self-debugging, only programs identified with negative labels during the iteration would perform further repair. Therefore, our focus is primarily on the distribution of different negative labels. We observed that testing on self-generated tests is more likely to result in false negative labels than true negative ones on easy problems. However, a different pattern emerges on medium and hard ones, where

false negatives outnumber true negatives<sup>2</sup>. *This discrepancy is primarily attributed to lower performance on more complex and challenging tasks, where negative labels from self-testing are more likely to align with the actual labels of the generated programs.* Relying solely on label feedback during self-debugging inadvertently reduces the bias from self-testing, thereby increasing the prevalence of true negative labels. However, when incorrect details are included in feedback, the performance declines compared to using only label feedback.

Moreover, when self-testing results in a false negative, it is crucial for the model to accurately identify the errors within the feedback and preserve the original answers. The efficacy of post-execution self-debugging, depends not only on the model’s ability to identify the defects in its own programs when presented with true negative labels but also on its ability to recognize the faulty execution feedback in the case of false negatives.

### 3.4 RQ3: In-Execution Reasoning Helps Self-Debugging

Despite the conditional utility of post-execution self-debugging in challenging tasks, its dependency on potentially unreliable test labels remains a fundamental fragility. We therefore explore in this subsection whether the in-execution method—by debugging through intermediate states—can achieve consistent gains across different levels of tasks.

Drawing inspiration from the implementation presented in LDB (Zhong et al., 2024), we divide an execution trace into several units based on basic blocks in the program’s control flow graph and collect the variables before and after these units during program execution to facilitate in-execution self-debugging. Please refer to Appendix E for a detailed illustration. Different from existing works that utilize runtime information for debugging (Zhong et al., 2024; Ni et al., 2024; Ding et al., 2024), we decouple post-execution signals from these methods, which means the labels and details from post-execution results are not accessible for the models. Therefore, the models must determine program correctness merely based on the test input and corresponding intermediate states.

### In-Execution Self-Debugging Shows Significant Potential for Improvement.

The results of in-

<sup>2</sup>We also present an analysis of test bias on HumanEval and MBPP in Figure 6.

Table 2: Results of **in-execution self-debugging** (in-exec.) with *self-generated tests*. The values in **red** are declines while those in **green** are increases compared to the initial generation (w/o self-debug).

Model	Method	#Iter	HumanEval+	MBPP+	LiveCodeBench			
					easy	medium	hard	overall
GPT-4o-2024-05-13	w/o self-debug	0	87.8	76.5	89.3	33.1	6.0	46.0
	in-exec. w/ trace	1	89.0 <sup>+1.2</sup>	77.8 <sup>+1.3</sup>	91.2 <sup>+1.9</sup>	34.9 <sup>+1.8</sup>	6.0 <sup>+0.0</sup>	47.3 <sup>+1.3</sup>
		2	88.4 <sup>+0.6</sup>	79.1 <sup>+2.6</sup>	91.8 <sup>+2.5</sup>	34.9 <sup>+1.8</sup>	6.0 <sup>+0.0</sup>	47.6 <sup>+1.6</sup>
Claude-3.5-Sonnet	w/o self-debug	0	89.0	77.0	93.1	48.0	16.4	55.8
	in-exec. w/ trace	1	89.6 <sup>+0.6</sup>	77.2 <sup>+0.2</sup>	95.0 <sup>+1.9</sup>	49.1 <sup>+1.1</sup>	17.2 <sup>+0.8</sup>	57.1 <sup>+1.3</sup>
		2	87.2 <sup>-1.8</sup>	76.2 <sup>-0.8</sup>	93.7 <sup>+0.6</sup>	48.6 <sup>+0.6</sup>	17.2 <sup>+0.8</sup>	56.4 <sup>+0.6</sup>
LLaMA-3-70B-Ins	w/o self-debug	0	73.8	71.2	72.3	10.3	2.6	30.2
	in-exec. w/ trace	1	70.1 <sup>-3.7</sup>	69.6 <sup>-1.6</sup>	73.0 <sup>+0.7</sup>	11.4 <sup>+1.1</sup>	3.4 <sup>+0.8</sup>	31.1 <sup>+0.9</sup>
		2	74.4 <sup>+0.6</sup>	69.6 <sup>-1.6</sup>	71.1 <sup>-1.2</sup>	12.0 <sup>+1.7</sup>	3.4 <sup>+0.8</sup>	30.7 <sup>+0.5</sup>
Qwen2.5-Coder-7B-Ins	w/o self-debug	0	81.7	70.6	74.8	23.4	8.6	35.8
	in-exec. w/ trace	1	82.3 <sup>+0.6</sup>	71.4 <sup>+0.8</sup>	75.5 <sup>+0.7</sup>	24.0 <sup>+0.6</sup>	8.6 <sup>+0.0</sup>	36.2 <sup>+0.4</sup>
		2	82.3 <sup>+0.6</sup>	72.0 <sup>+1.4</sup>	76.1 <sup>+1.3</sup>	24.0 <sup>+0.6</sup>	8.6 <sup>+0.0</sup>	36.4 <sup>+0.6</sup>

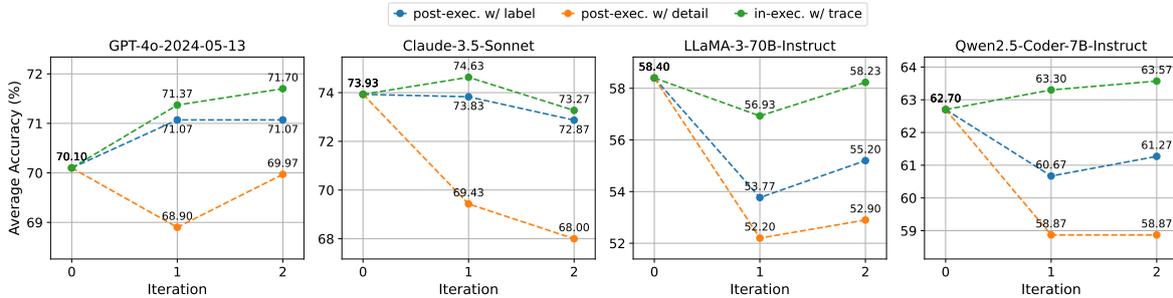


Figure 4: Average accuracy of both post-execution and in-execution self-debugging across three benchmarks over two iterations. In-execution self-debugging, which only leverages intermediate runtime information, consistently outperforms the post-execution approach.

439 execution self-debugging are detailed in Table 2.  
 440 For basic problems from HumanEval and MBPP,  
 441 we observe that self-debugging gains notable im-  
 442 provement for GPT-4o and Qwen2.5-Coder-7B-  
 443 Instruct when utilizing in-execution information.  
 444 Claude-3.5-Sonnet also shows performance im-  
 445 provements at the first iteration on both bench-  
 446 marks. For LLaMA-3-70B-Instruct, the perfor-  
 447 mance on HumanEval surpasses the baseline at the  
 448 second iteration. However, there is a slight degrada-  
 449 tion in performance for certain tasks and iterations  
 450 compared to the initial generation. For competi-  
 451 tive problems from LiveCodeBench, the results  
 452 demonstrate the efficacy of the in-execution self-  
 453 debugging across all models. The results indicate  
 454 that in-execution self-debugging is a potentially  
 455 effective approach for improving performance on  
 456 diverse programming problems.

#### 457 How In-Execution Works Compared to Post- 458 Execution Method

We present the average accu-

racy of both post-execution and in-execution meth-  
 ods over two iterations in Figure 4. Since the debug-  
 ging process is triggered using self-generated tests  
 rather than ground truth oracle tests, performance  
 does not necessarily improve with successive itera-  
 tions. This differs from other existing debugging  
 methods (Chen et al., 2024b; Shinn et al., 2023;  
 Zhong et al., 2024), where the set of positive pro-  
 grams in the current iteration is typically a superset  
 of those in the previous iteration. The results sug-  
 gest that *in-execution self-debugging, which only  
 leverages intermediate runtime information, con-  
 sistently outperforms the post-execution approach.*  
*Moreover, in most cases, it shows improvements  
 compared to the initial generation, indicating sig-  
 nificant potential for future development.* To further  
 illustrate how the in-execution method works, we  
 provide a comparison with an example from Hu-  
 manEval using GPT-4o in Appendix F.

459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477

## Strengths and Weaknesses of In-Execution Self-Debugging

To conclude, post-execution self-debugging utilizes final execution results to reflect upon and debug programs. However, the unreliability of the self-generated tests could bias the model away from the correct answer. *Although this can provide some relief on challenging tasks, it is not a long-term solution, especially when those hard tasks can also be solved well over time.* On the contrary, in-execution self-debugging allows the models to perform fine-grained feedback solely on the intermediate states during the execution process, mitigating the test bias introduced by potentially flawed tests. It shows the potential to better align the programs with the requirements in real-world scenarios. Nonetheless, self-debugging with in-execution information depends heavily on the LLMs’ code reasoning capabilities and lacks formal guarantees of success, as the pass rate drops for LLaMA-3-70B-Instruct on MBPP. We expect that improvements in LLM capabilities will enhance the efficacy of this paradigm. Based on our findings, we deliver a detailed discussion on directions for future work in Appendix A.

## 4 Related Work

**LLM-Based Code Generation.** Code generation is the automatic production of source code based on natural language descriptions. Researchers have proposed various approaches to enhance the quality of code generated by these models. Some works, like LLaMA (Touvron et al., 2023a,b; Dubey et al., 2024) and DeepSeek (DeepSeek-AI et al., 2024b,a) series, focus on optimizing model training, while others aim to improve code quality through post-processing techniques (Chen et al., 2023; Zhang et al., 2023b; Inala et al., 2022). Among these post-processing techniques, methods that involve self-debugging have gained considerable attention. Through feedback from execution results, self-debugging allows models to autonomously debug and refine previously generated code, enhancing the final output. Self-debugging has been integrated into various LLM-based code generation methods (Yang et al., 2024; Zhang et al., 2024; Dong et al., 2023; Huang et al., 2023). In this work, we revisit these techniques, evaluating self-debugging with self-generated tests on both basic and competitive tasks.

**Self-Debugging with LLMs.** As LLMs have evolved, the idea of using models to refine their

own output has become more popular. Most of these methods rely on prompting LLMs with execution results to improve the code. These methods often rely on pre-existing or generated tests to execute the code, capturing execution information that is then used to refine the output code ((Olausson et al., 2024; Wang et al., 2024; Dong et al., 2023; Madaan et al., 2023; Zhang et al., 2023a; Jiang et al., 2024)). Self-Refine (Madaan et al., 2023) conducts a broad evaluation of self-debugging in code models, highlighting that performance can be improved with higher-quality feedback or human intervention. Self-Debugging (Chen et al., 2024b) and Self-Edit (Zhang et al., 2023a) use the framework in which LLMs can iteratively debug their own generated code by utilizing execution results. LDB (Zhong et al., 2024) utilizes runtime execution information to help debug generated programs. To enhance the execution understanding capability of LLMs, SemCoder (Ding et al., 2024) trains models to simulate execution, enabling abstract reasoning and debugging. NExT (Ni et al., 2024) trains LLMs to generate NL rationales based on execution traces for accurate runtime reasoning. These methods leverage both post and in-execution information for self-debugging since there are pre-defined oracle tests, which is an idealized setup in real-world scenarios. In this work, we aim to investigate the effect of post-execution and in-execution information respectively, particularly with self-generated tests. We propose a unified framework and provide a detailed analysis of these methods.

## 5 Conclusion

In this work, we explore self-debugging with self-generated tests to enhance the code generation capabilities of large language models (LLMs). We introduce and evaluate two paradigms: post-execution and in-execution self-debugging. While post-execution methods struggle with simpler tasks due to bias from self-generated tests, in-execution self-debugging leverages intermediate runtime states to mitigate this bias and consistently outperforms post-execution approaches. Our findings highlight the importance of reducing reliance on flawed self-generated tests and incorporating richer runtime feedback for more accurate and robust program refinement. This work underscores the potential of in-execution self-debugging to advance LLMs in code generation, paving the way for more autonomous and reliable code generation systems.

## 578 Limitations

579 We evaluate the efficacy of the proposed self-  
580 debugging paradigms on a wide range of LLMs and  
581 across three benchmarks covering both basic and  
582 competitive programming problems. While our ex-  
583 periments could provide valuable insights into self-  
584 debugging for real-world applications, it remains  
585 unclear whether these findings would directly trans-  
586 late to more complex software engineering scenar-  
587 ios. For instance, consider repository-level debug-  
588 ging: a model analyzing intermediate API call se-  
589 quences or variable propagation across modules  
590 could mirror the in-execution paradigm. However,  
591 questions such as how to efficiently capture and  
592 structure runtime data across interconnected mod-  
593 ules, or how to handle dependencies and interac-  
594 tions between components, remain open. These  
595 complexities highlight the need for further research  
596 to adapt self-debugging techniques to large-scale,  
597 real-world software systems.

598 We leave these challenges for future work, as  
599 they represent an important next step in advanc-  
600 ing the applicability of self-debugging paradigms  
601 beyond self-contained programming tasks. By ad-  
602 dressing these issues, we can move closer to de-  
603 veloping autonomous and self-evolving systems  
604 capable of handling the intricacies of modern soft-  
605 ware engineering.

## 606 References

- 607 Jacob Austin, Augustus Odena, Maxwell I. Nye,  
608 Maarten Bosma, Henryk Michalewski, David Dohan,  
609 Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le,  
610 and Charles Sutton. 2021. [Program synthesis with  
611 large language models](#). *CoRR*, abs/2108.07732.
- 612 Islem Bouzenia, Yangruibo Ding, Kexin Pei, Baishakhi  
613 Ray, and Michael Pradel. 2023. [Tracefixer:  
614 Execution trace-driven program repair](#). *CoRR*,  
615 abs/2304.12743.
- 616 Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner,  
617 Bowen Baker, Leo Gao, Leopold Aschenbrenner,  
618 Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan  
619 Leike, Ilya Sutskever, and Jeffrey Wu. 2024. [Weak-  
620 to-strong generalization: Eliciting strong capabilities  
621 with weak supervision](#). In *Forty-first International  
622 Conference on Machine Learning, ICML 2024, Vi-  
623 enna, Austria, July 21-27, 2024*. OpenReview.net.
- 624 Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan,  
625 Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2023.  
626 [Codet: Code generation with generated tests](#). In  
627 *The Eleventh International Conference on Learning  
628 Representations, ICLR 2023, Kigali, Rwanda, May  
629 1-5, 2023*. OpenReview.net.

- Junkai Chen, Zhiyuan Pan, Xing Hu, Zhenhao Li, Ge Li,  
and Xin Xia. 2024a. [Evaluating large language  
models with runtime behavior of program execution](#).  
*CoRR*, abs/2403.16437. 630  
631  
632  
633
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming  
Yuan, Henrique Pondé de Oliveira Pinto, Jared Ka-  
plan, Harri Edwards, Yuri Burda, Nicholas Joseph,  
Greg Brockman, Alex Ray, Raul Puri, Gretchen  
Krueger, Michael Petrov, Heidy Khlaaf, Girish Sas-  
try, Pamela Mishkin, Brooke Chan, Scott Gray,  
Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz  
Kaiser, Mohammad Bavarian, Clemens Winter,  
Philippe Tillet, Felipe Petroski Such, Dave Cum-  
mings, Matthias Plappert, Fotios Chantzis, Eliza-  
beth Barnes, Ariel Herbert-Voss, William Hebg-  
Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie  
Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain,  
William Saunders, Christopher Hesse, Andrew N.  
Carr, Jan Leike, Joshua 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. [Evaluat-  
ing large language models trained on code](#). *CoRR*,  
abs/2107.03374. 634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654
- Xinyun Chen, Maxwell Lin, Nathanael Schärli, and  
Denny Zhou. 2024b. [Teaching large language mod-  
els to self-debug](#). In *The Twelfth International Con-  
ference on Learning Representations, ICLR 2024,  
Vienna, Austria, May 7-11, 2024*. OpenReview.net. 655  
656  
657  
658  
659
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingx-  
uan Wang, Bochao Wu, Chengda Lu, Chenggang  
Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,  
Damai Dai, Daya Guo, Dejian Yang, Deli Chen,  
Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai,  
Fuli Luo, Guangbo Hao, Guanting Chen, Guowei  
Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng  
Wang, Haowei Zhang, Honghui Ding, Huajian Xin,  
Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang,  
Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang,  
Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie  
Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu,  
Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean  
Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao,  
Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang,  
Mingchuan Zhang, Minghua Zhang, Minghui Tang,  
Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang,  
Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu  
Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge,  
Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin  
Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao  
Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu,  
Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu  
Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou,  
Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun,  
W. L. Xiao, and Wangding Zeng. 2024a. [Deepseek-  
v3 technical report](#). *CoRR*, abs/2412.19437. 660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686
- DeepSeek-AI, Qihao Zhu, Daya Guo, Zhihong Shao,  
Dejian Yang, Peiyi Wang, Runxin Xu, Y. Wu, Yukun  
Li, Huazuo Gao, Shirong Ma, Wangding Zeng, Xiao  
Bi, Zihui Gu, Hanwei Xu, Damai Dai, Kai Dong, 687  
688  
689  
690

691	Liyue Zhang, Yishi Piao, Zhibin Gou, Zhenda Xie,	<i>ICML 2024, Vienna, Austria, July 21-27, 2024</i> . Open-	751
692	Zhewen Hao, Bingxuan Wang, Junxiao Song, Deli	Review.net.	752
693	Chen, Xin Xie, Kang Guan, Yuxiang You, Aixin		
694	Liu, Qiushi Du, Wenjun Gao, Xuan Lu, Qinyu Chen,	Dong Huang, Qingwen Bu, Jie M. Zhang, Michael Luck,	753
695	Yaohui Wang, Chengqi Deng, Jiashi Li, Chenggang	and Heming Cui. 2023. <a href="#">Agentcoder: Multi-agent-</a>	754
696	Zhao, Chong Ruan, Fuli Luo, and Wenfeng Liang.	<a href="#">based code generation with iterative testing and opti-</a>	755
697	2024b. <a href="#">Deepseek-coder-v2: Breaking the barrier of</a>	<a href="#">misation</a> . <i>CoRR</i> , abs/2312.13010.	756
698	<a href="#">closed-source models in code intelligence</a> . <i>CoRR</i> ,		
699	abs/2406.11931.		
700	Yangruibo Ding, Jinjun Peng, Marcus J. Min, Gail E.	Jie Huang, Xinyun Chen, Swaroop Mishra,	757
701	Kaiser, Junfeng Yang, and Baishakhi Ray. 2024.	Huaixiu Steven Zheng, Adams Wei Yu, Xiny-	758
702	<a href="#">Semcoder: Training code language models with com-</a>	ing Song, and Denny Zhou. 2024. <a href="#">Large language</a>	759
703	<a href="#">prehensive semantics reasoning</a> . In <i>Advances in Neu-</i>	<a href="#">models cannot self-correct reasoning yet</a> . In <i>The</i>	760
704	<i>ral Information Processing Systems 38: Annual Con-</i>	<i>Twelfth International Conference on Learning</i>	761
705	<i>ference on Neural Information Processing Systems</i>	<i>Representations, ICLR 2024, Vienna, Austria, May</i>	762
706	<i>2024, NeurIPS 2024, Vancouver, BC, Canada, De-</i>	<i>7-11, 2024</i> . OpenReview.net.	763
707	<i>ember 10 - 15, 2024</i> .		
708	Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. 2023. <a href="#">Self-</a>	Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Day-	764
709	<a href="#">collaboration code generation via chatgpt</a> . <i>CoRR</i> ,	iheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang,	765
710	abs/2304.07590.	Bowen Yu, Kai Dang, An Yang, Rui Men, Fei Huang,	766
		Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and	767
		Junyang Lin. 2024. <a href="#">Qwen2.5-coder technical report</a> .	768
		<i>CoRR</i> , abs/2409.12186.	769
711	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,	Jeevana Priya Inala, Chenglong Wang, Mei Yang, An-	770
712	Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,	drés Codas, Mark Encarnación, Shuvendu K. Lahiri,	771
713	Akhil Mathur, Alan Schelten, Amy Yang, Angela	Madanlal Musuvathi, and Jianfeng Gao. 2022. <a href="#">Fault-</a>	772
714	Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang,	<a href="#">aware neural code rankers</a> . In <i>Advances in Neural</i>	773
715	Archi Mitra, Archie Sravankumar, Artem Korenev,	<i>Information Processing Systems 35: Annual Confer-</i>	774
716	Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien	<i>ence on Neural Information Processing Systems 2022,</i>	775
717	Rodriguez, Austen Gregerson, Ava Spataru, Bap-	<i>NeurIPS 2022, New Orleans, LA, USA, November 28</i>	776
718	tiste Rozière, Bethany Biron, Binh Tang, Bobbie	<i>- December 9, 2022</i> .	777
719	Chern, Charlotte Caucheteux, Chaya Nayak, Chloe		
720	Bi, Chris Marra, Chris McConnell, Christian Keller,	Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia	778
721	Christophe Touret, Chunyang Wu, Corinne Wong,	Yan, Tianjun Zhang, Sida Wang, Armando Solar-	779
722	Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-	Lezama, Koushik Sen, and Ion Stoica. 2024. <a href="#">Live-</a>	780
723	lonsius, Daniel Song, Danielle Pintz, Danny Livshits,	<a href="#">codebench: Holistic and contamination free eval-</a>	781
724	David Esiobu, Dhruv Choudhary, Dhruv Mahajan,	<a href="#">uation of large language models for code</a> . <i>CoRR</i> ,	782
725	Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,	abs/2403.07974.	783
726	Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,		
727	Emily Dinan, Eric Michael Smith, Filip Radenovic,	Nan Jiang, Xiaopeng Li, Shiqi Wang, Qiang Zhou,	784
728	Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Geor-	Soneya Binta Hossain, Baishakhi Ray, Varun Ku-	785
729	gia Lewis Anderson, Graeme Nail, Grégoire Mialon,	mar, Xiaofei Ma, and Anoop Deoras. 2024. <a href="#">Training</a>	786
730	Guan Pang, Guillem Cucurell, Hailey Nguyen, Han-	<a href="#">llms to better self-debug and explain code</a> . <i>CoRR</i> ,	787
731	nah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov,	abs/2405.18649.	788
732	Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan		
733	Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan	Shuyang Jiang, Yuhao Wang, and Yu Wang. 2023. <a href="#">Self-</a>	789
734	Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,	<a href="#">evolve: A code evolution framework via large lan-</a>	790
735	Jeet Shah, Jelmer van der Linde, Jennifer Billock,	<a href="#">guage models</a> . <i>CoRR</i> , abs/2306.02907.	791
736	Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi,		
737	Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu,	Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Ling-	792
738	Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph	ming Zhang. 2023. <a href="#">Is your code generated by chatgpt</a>	793
739	Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia,	<a href="#">really correct? rigorous evaluation of large language</a>	794
740	Kalyan Vasuden Alwala, Kartikeya Upasani, Kate	<a href="#">models for code generation</a> . In <i>Advances in Neural</i>	795
741	Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and	<i>Information Processing Systems 36: Annual Confer-</i>	796
742	et al. 2024. <a href="#">The llama 3 herd of models</a> . <i>CoRR</i> ,	<i>ence on Neural Information Processing Systems 2023,</i>	797
743	abs/2407.21783.	<i>NeurIPS 2023, New Orleans, LA, USA, December 10</i>	798
		<i>- 16, 2023</i> .	799
744	Robert Gold. 2010. <a href="#">Control flow graphs and code cover-</a>	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	800
745	<a href="#">age</a> . <i>Int. J. Appl. Math. Comput. Sci.</i> , 20(4):739–749.	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	801
		Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	802
746	Alex Gu, Baptiste Rozière, Hugh James Leather, Ar-	Shashank Gupta, Bodhisattwa Prasad Majumder,	803
747	mando Solar-Lezama, Gabriel Synnaeve, and Sida	Katherine Hermann, Sean Welleck, Amir Yazdan-	804
748	Wang. 2024. <a href="#">Cruxeval: A benchmark for code</a>	bakhsh, and Peter Clark. 2023. <a href="#">Self-refine: Itera-</a>	805
749	<a href="#">reasoning, understanding and execution</a> . In <i>Forty-</i>	<a href="#">tive refinement with self-feedback</a> . In <i>Advances in</i>	806
750	<i>first International Conference on Machine Learning</i> ,	<i>Neural Information Processing Systems 36: Annual</i>	807

808					
809					
810					
811	Ansong Ni, Miltiadis Allamanis, Arman Cohan, Yinlin				
812	Deng, Kensen Shi, Charles Sutton, and Pengcheng				
813	Yin. 2024. <a href="#">Next: Teaching large language models to</a>				
814	<a href="#">reason about code execution</a> . In <i>Forty-first Interna-</i>				
815	<i>tional Conference on Machine Learning, ICML 2024,</i>				
816	<i>Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.				
817	Maxwell I. Nye, Anders Johan Andreassen, Guy Gur-				
818	Ari, Henryk Michalewski, Jacob Austin, David				
819	Bieber, David Dohan, Aitor Lewkowycz, Maarten				
820	Bosma, David Luan, Charles Sutton, and Augustus				
821	Odena. 2021. <a href="#">Show your work: Scratchpads for inter-</a>				
822	<a href="#">mediate computation with language models</a> . <i>CoRR</i> ,				
823	abs/2112.00114.				
824	Theo X. Olausson, Jeevana Priya Inala, Chenglong				
825	Wang, Jianfeng Gao, and Armando Solar-Lezama.				
826	2024. <a href="#">Is self-repair a silver bullet for code genera-</a>				
827	<a href="#">tion?</a> In <i>The Twelfth International Conference on</i>				
828	<i>Learning Representations, ICLR 2024, Vienna, Aus-</i>				
829	<i>tria, May 7-11, 2024</i> . OpenReview.net.				
830	Tal Ridnik, Dedy Kredo, and Itamar Friedman.				
831	2024. <a href="#">Code generation with alphacodium: From</a>				
832	<a href="#">prompt engineering to flow engineering</a> . <i>CoRR</i> ,				
833	abs/2401.08500.				
834	Noah Shinn, Federico Cassano, Ashwin Gopinath,				
835	Karthik Narasimhan, and Shunyu Yao. 2023. <a href="#">Re-</a>				
836	<a href="#">flexion: language agents with verbal reinforcement</a>				
837	<a href="#">learning</a> . In <i>Advances in Neural Information Pro-</i>				
838	<i>cessing Systems 36: Annual Conference on Neural</i>				
839	<i>Information Processing Systems 2023, NeurIPS 2023,</i>				
840	<i>New Orleans, LA, USA, December 10 - 16, 2023</i> .				
841	Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu				
842	Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei Huang,				
843	Dacheng Tao, and Jingren Zhou. 2024. <a href="#">A survey</a>				
844	<a href="#">on self-evolution of large language models</a> . <i>CoRR</i> ,				
845	abs/2404.14387.				
846	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier				
847	Martinet, Marie-Anne Lachaux, Timothée Lacroix,				
848	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal				
849	Azhar, Aurélien Rodriguez, Armand Joulin, Edouard				
850	Grave, and Guillaume Lample. 2023a. <a href="#">Llama: Open</a>				
851	<a href="#">and efficient foundation language models</a> . <i>CoRR</i> ,				
852	abs/2302.13971.				
853	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-				
854	bert, Amjad Almahairi, Yasmine Babaei, Nikolay				
855	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti				
856	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-				
857	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,				
858	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,				
859	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-				
860	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan				
861	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,				
862	Isabel Kloumann, Artem Korenev, Punit Singh Koura,				
	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-				
	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-				
	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-				
	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-				
	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,				
	Ruan Silva, Eric Michael Smith, Ranjan Subrama-				
	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-				
	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,				
	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,				
	Melanie Kambadur, Sharan Narang, Aurélien Ro-				
	driguez, Robert Stojnic, Sergey Edunov, and Thomas				
	Scialom. 2023b. <a href="#">Llama 2: Open foundation and</a>				
	<a href="#">fine-tuned chat models</a> . <i>CoRR</i> , abs/2307.09288.				
	Hanbin Wang, Zhenghao Liu, Shuo Wang, Ganqu Cui,				
	Ning Ding, Zhiyuan Liu, and Ge Yu. 2024. <a href="#">INTER-</a>				
	<a href="#">VENOR: prompting the coding ability of large lan-</a>				
	<a href="#">guage models with the interactive chain of repair</a> . In				
	<i>Findings of the Association for Computational Lin-</i>				
	<i>guistics, ACL 2024, Bangkok, Thailand and virtual</i>				
	<i>meeting, August 11-16, 2024</i> , pages 2081–2107.				
	John Yang, Carlos E. Jimenez, Alexander Wettig, Kil-				
	ian Lieret, Shunyu Yao, Karthik Narasimhan, and				
	Ofir Press. 2024. <a href="#">Swe-agent: Agent-computer inter-</a>				
	<a href="#">faces enable automated software engineering</a> . <i>CoRR</i> ,				
	abs/2405.15793.				
	Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin.				
	2024. <a href="#">Codeagent: Enhancing code generation with</a>				
	<a href="#">tool-integrated agent systems for real-world repo-</a>				
	<a href="#">level coding challenges</a> . In <i>Proceedings of the 62nd</i>				
	<i>Annual Meeting of the Association for Computa-</i>				
	<i>tional Linguistics (Volume 1: Long Papers), ACL</i>				
	<i>2024, Bangkok, Thailand, August 11-16, 2024</i> , pages				
	13643–13658. Association for Computational Lin-				
	guistics.				
	Kechi Zhang, Zhuo Li, Jia Li, Ge Li, and Zhi Jin. 2023a.				
	<a href="#">Self-edit: Fault-aware code editor for code genera-</a>				
	<a href="#">tion</a> . In <i>Proceedings of the 61st Annual Meeting of</i>				
	<i>the Association for Computational Linguistics (Vol-</i>				
	<i>ume 1: Long Papers), ACL 2023, Toronto, Canada,</i>				
	<i>July 9-14, 2023</i> , pages 769–787.				
	Tianyi Zhang, Tao Yu, Tatsunori Hashimoto, Mike				
	Lewis, Wen-Tau Yih, Daniel Fried, and Sida Wang.				
	2023b. <a href="#">Coder reviewer reranking for code genera-</a>				
	<a href="#">tion</a> . In <i>International Conference on Machine Learn-</i>				
	<i>ing, ICML 2023, 23-29 July 2023, Honolulu, Hawaii,</i>				
	<i>USA</i> , pages 41832–41846.				
	Li Zhong, Zilong Wang, and Jingbo Shang. 2024. <a href="#">De-</a>				
	<a href="#">bug like a human: A large language model debugger</a>				
	<a href="#">via verifying runtime execution step by step</a> . In <i>Find-</i>				
	<i>ings of the Association for Computational Linguistics,</i>				
	<i>ACL 2024, Bangkok, Thailand and virtual meeting,</i>				
	<i>August 11-16, 2024</i> , pages 851–870. Association for				
	Computational Linguistics.				

## A Directions for Future Work

**High-Quality Test Generation.** In this work, we demonstrate that post-execution self-debugging with self-generated tests struggles on basic problems due to biased evaluations, despite the significant potential shown by LLMs in automated test generation. This highlights the necessity for the research community to focus on the quality of LLM-generated tests before utilizing execution feedback derived from them. Developing techniques that enhance high-quality test synthesis is crucial to mitigate bias for post-execution self-debugging. It could be beneficial to implement an iterative refinement process wherein execution information is leveraged to improve the tests. This could involve using techniques like test-driven development where tests are continuously updated based on code changes and debugging outcomes.

**More Deliberate Runtime Information.** As demonstrated in Section 3.4, leveraging enriched runtime information from execution is a promising avenue for self-debugging. In particular, in-execution self-debugging has shown superior performance compared to post-execution in certain tasks, suggesting that more nuanced and reliable feedback leads to better performance. Designing more sophisticated methods for collecting and analyzing runtime information is a promising direction for further enhancing self-debugging capabilities. For instance, improving the intelligibility of execution trace representations for LLMs may prove beneficial (Ni et al., 2024). Additionally, beyond variables, other types of runtime information, such as code coverage, execution paths, and memory usage, could be effectively utilized to provide a more comprehensive view of program behavior (Chen et al., 2024a).

**Training for More Advanced Capabilities of LLMs.** Effective self-debugging with self-generated tests hinges on several core capabilities of LLMs. In terms of refinement, the model should be capable of accurately recognizing and localizing faults within the program. This requires not only an understanding of the code’s syntax but also its semantic logic and runtime behavior. Additionally, more advanced reasoning capabilities are needed to thoroughly analyze execution feedback. The model should comprehend the relationship between the code logic and the feedback, thereby deducing the runtime structure of program statements and vari-

ables. Training LLMs with diverse debugging scenarios (Ni et al., 2024) and incorporating multi-step reasoning frameworks could further enhance their ability to diagnose and fix errors autonomously.

**Self-Evolution of Code LLMs.** Self-debugging opens up possibilities for developing more advanced LLMs without reliance on human supervision or guidance from stronger models (Burns et al., 2024). Traditionally, human-generated test cases serve as a strong supervisory signal for aligning code generation, but the collection of such tests is labor-intensive, leading to a sparsity of labeled data for effective code refinement. Self-generated tests, by contrast, offer a viable path for self-improvement (Tao et al., 2024). They alleviate the burden of manual test generation and pave the way toward truly autonomous self-correcting code generation systems (Chen et al., 2024b). By iteratively refining both the code and the tests, LLMs can evolve into more robust and reliable systems capable of handling increasingly complex programming tasks. This self-evolution paradigm could revolutionize the development of code generation models, making them more adaptable and scalable for real-world applications.

These directions highlight the potential for advancing self-debugging techniques, ultimately leading to more autonomous, accurate, and efficient code generation systems. By addressing these challenges, the research community can unlock new possibilities for LLMs in software development and beyond.

## B Benchmark

We utilize HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021) and LiveCodeBench (Jain et al., 2024) covering basic and competitive programming problems to comprehensively evaluate the efficacy of self-debugging:

- **HumanEval and MBPP** HumanEval (Chen et al., 2021) consists of 164 programming problems written by humans. Each problem provides a Python function signature and a docstring as its specification. MBPP (Austin et al., 2021) includes 974 programming problems written by contributors through crowdsourcing. Each of these problems features a problem statement, a function signature, and three example tests. To enhance the reliability and accuracy of evaluations, EvalPlus (Liu et al., 2023) extends HumanEval into a more comprehensive version

known as HumanEval+ with 80 times more tests than the original HumanEval. Similarly, MBPP+ is an augmentation of the original MBPP, offering 35 times more tests. In our experiments, we use the latest version of MBPP+, which consists of 378 programming problems.

- **LiveCodeBench** LiveCodeBench (Jain et al., 2024) is a contamination-free benchmark that continuously collects new problems from prominent competitive programming platforms. As of now, LiveCodeBench features a collection of over 600 high-quality programming problems. These problems encompass a wide range of difficulty levels and topics, providing a comprehensive evaluation for the coding capabilities of LLMs. In our experiments, we select 450 problems that were published between September 2023 and September 2024.

## C Validation of Generated Tests

In Section 3.2, we leverage program contracts and canonical solutions to evaluate the reliability of the generated tests for problems in HumanEval and MBPP<sup>3</sup>. Program contracts ensure that the generated test inputs adhere to the preconditions specified for each problem, while canonical solutions serve as the ground truth to verify the correctness of the test outputs. We present examples for both benchmarks in Figure 5.

To validate test inputs, we place contracts at the beginning of the function and pass the test input to it. If no assertion error occurs, the test input is considered valid. Some functions have additional requirements for inputs beyond just type, so the contracts check for both type mismatches and any other specified constraints. For test output validation, we use canonical solutions to collect the execution output for a given valid input and compare it with the generated test output to confirm correctness.

## D Varying Number of Generated Tests

To investigate the effect of the number of self-generated tests, we employ GPT-4o-2024-05-13 to generate varying number of tests  $N = [10, 15, 20]$  for each programming problem in HumanEval and MBPP. Following Section 3.2, we calculate the accuracy of these generated tests and the results

<sup>3</sup>Neither program contracts nor canonical solutions are provided by LiveCodeBench.

Table 3: Results for increasing number of self-generated tests using GPT-4o-2024-05-13 on HumanEval and MBPP.

(a) HumanEval			
#Num of Tests	Input	Output	Suite
10	97.63%	89.77%	59.15%
15	97.89%	88.86%	52.44% (-6.71%)
20	98.11%	86.01%	48.17% (-10.98%)
(b) MBPP			
#Num of Tests	Input	Output	Suite
10	94.81%	85.60%	58.73%
15	94.96%	85.27%	53.70% (-5.03%)
20	95.10%	82.94%	50.53% (-8.20%)

Table 4: Results of post-execution self-debugging with detailed feedback and in-execution self-debugging on HumanEval and MBPP when using different sizes of the self-generated test suite. The values highlighted in red or green are changes relative to the initial generation.

Method	#N	HumanEval+	MBPP+
w/o self-debug	0	87.8	76.5
post-exec. w/ detail	10	84.1 <sup>-3.7</sup>	76.2 <sup>-0.3</sup>
	15	84.1 <sup>-3.7</sup>	75.9 <sup>-0.6</sup>
	20	83.5 <sup>-4.3</sup>	75.9 <sup>-0.6</sup>
in-exec. w/ trace	10	89.0 <sup>+1.2</sup>	77.8 <sup>+1.3</sup>
	15	88.4 <sup>+0.6</sup>	78.0 <sup>+1.5</sup>
	20	88.4 <sup>+0.6</sup>	77.2 <sup>+0.7</sup>

are summarized in Table 3. As the number of self-generated tests increases, the presence of more challenging edge cases also rises, leading to a reduction in the overall accuracy of the test suites. Specifically, when the model generates up to 20 tests per problem, the accuracy of the test suite decreases from 59.15% to 48.17% for HumanEval and from 58.73% to 50.53% for MBPP.

We further evaluate the performance of both post-execution self-debugging with detailed feedback and in-execution self-debugging. The results in Table 4 indicate that with an increased number of self-generated tests, the performance of post-execution self-debugging experiences a slight decline on both HumanEval and MBPP. This decline is attributed to the lower accuracy of the test suites, which increases the rate of false negatives and thereby hinders the efficacy of post-execution self-debugging. In contrast, in-execution self-debugging, which leverages intermediate runtime traces, demonstrates consistent improvement over the initial generation.

These findings highlight the robustness of in-

```

# HumanEval/21

def rescale_to_unit(numbers: List[float]) -> List[float]:
    """ Given list of numbers (of at least two elements), apply a linear transform to
    that list, such that the smallest number will become 0 and the largest will become 1
    >>> rescale_to_unit([1.0, 2.0, 3.0, 4.0, 5.0])
    [0.0, 0.25, 0.5, 0.75, 1.0]
    """

    assert all(type(x) in [int, float] for x in numbers), "invalid inputs" #
    $_CONTRACT_$
    assert len(numbers) >= 2, "invalid inputs" # $_CONTRACT_$
    assert max(numbers) > min(numbers), "invalid inputs" # $_CONTRACT_$

    ma, mi = max(numbers), min(numbers)
    k = 1 / (ma - mi)
    return list(map(lambda x: (x - mi) * k, numbers))

```

```

# MBPP/439

"""
Write a function to join a list of multiple integers into a single integer.
assert multiple_to_single([11, 33, 50])==113350
"""

def multiple_to_single(L):

    assert isinstance(L, list), "invalid inputs" # $_CONTRACT_$
    assert len(L) > 0, "invalid inputs" # $_CONTRACT_$
    assert all(isinstance(item, int) for item in L), "invalid inputs" # $_CONTRACT_$
    assert all(item > 0 for item in L[1:]), "invalid inputs" # $_CONTRACT_$

    return int(''.join(map(str,L)))

```

Figure 5: Examples of program contracts in HumanEval and MBPP. Program contracts consist of assertions that specify conditions necessary for a valid input.

1085 execution self-debugging in handling the chal- 1104  
1086 lenges posed by less accurate test suites, as it relies 1105  
1087 on precise runtime information rather than poten- 1106  
1088 tially flawed self-generated tests. They also under- 1107  
1089 score the importance of balancing the quantity and 1108  
1090 quality of self-generated tests to ensure their ef- 1109  
1091 fectiveness in evaluating and debugging programs. 1110  
1092 By leveraging runtime insights, in-execution self- 1111  
1093 debugging provides a more reliable approach to 1112  
1094 program repair and enhancement, particularly in 1113  
1095 scenarios where test accuracy is compromised. 1114

## 1096 E Collection of In-Execution Information 1115

1097 In this work, we partition a program into basic 1116  
1098 blocks based on the nodes within its control flow 1117  
1099 graph (CFG) (Gold, 2010; Zhong et al., 2024). A 1118  
1100 control flow graph is a representation of the differ- 1119  
1101 ent blocks of code in a program and various execu- 1120  
1102 tion paths that an executor can take through the 1121  
1103 code. Each block represents a maximal sequence 1122  
1123

of statements and expressions without any jumps 1104  
or jump targets. Directed edges in the graph reflect 1105  
the possible execution order among the blocks. 1106

We simply use variables as the key indicators 1107  
of intermediate states during the execution. To 1108  
facilitate the visibility of necessary variables, we 1109  
strategically insert specific lines of code before and 1110  
after these basic blocks. Given a test input, the 1111  
program subsequently outputs the variables in ex- 1112  
ecution order. We then organize the content of 1113  
each reachable basic block and its corresponding 1114  
variables into execution trace blocks in sequence, 1115  
which are integrated into the debugging prompt 1116  
as shown in Figure 10. We provide an example 1117  
illustrating the representation of in-execution infor- 1118  
mation in Figure 7. For programs containing loops, 1119  
we are able to capture the change of the intermedi- 1120  
ate states within each loop iteration. Additionally, 1121  
in cases where the number of trace blocks exceeds 1122  
a predefined threshold  $N$ , we adopt a strategy to 1123

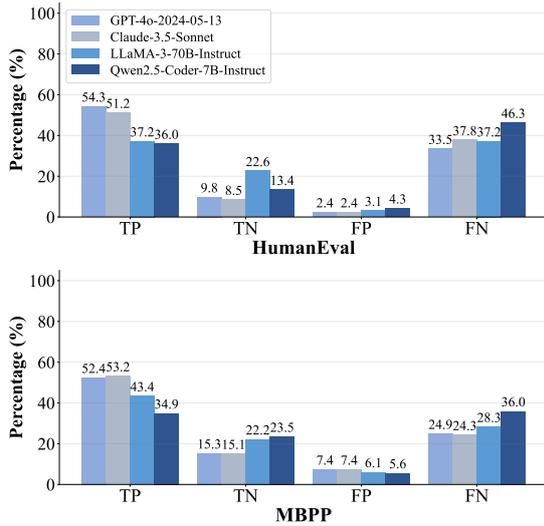


Figure 6: Test bias in labels when evaluating the programs with self-generated tests for HumanEval and MBPP.

retain the first  $N/2$  blocks and the last  $N/2$  blocks of the trace. It ensures that both the initial and final segments of the trace, which may contain essential information, are preserved for analysis while maintaining a manageable context length.

## F Case Study

In our experiments, we observe that in-execution self-debugging, which leverages runtime information, consistently outperforms post-execution one across various levels of programming tasks. To better understand the unique strengths and weaknesses of these two paradigms, we provide an example involving GPT-4o in Figure 7.

This example illustrates different outcomes of post-execution self-debugging with detailed test feedback and in-execution self-debugging with execution traces. The initial completion for the `can_arrange` function is correct. However, it is evaluated against an erroneous self-generated test that, according to the specification, should return 4 instead of 1. This discrepancy causes the model to misinterpret the condition in the problem, thereby leading to an incorrect revision of the program. Feedback from post-execution self-debugging, when based on flawed self-generated tests, biases the model away from the problem’s specification, resulting in a faulty program.

In contrast, in-execution self-debugging leverages test inputs and their corresponding runtime information to assess program correctness. As depicted in Figure 7, this approach enables the model

to perform a fine-grained analysis on the execution trace block by block without being influenced by the potential bias introduced by self-generated tests. By examining the intermediate states in execution order, the model confirms that the trace aligns with the expected behavior of the function, ultimately preserving the correctness of the original program.

This example highlights the critical advantage of in-execution self-debugging: its ability to rely on precise runtime information rather than potentially flawed self-generated tests. By avoiding the test bias, in-execution self-debugging provides a more reliable and effective approach to program repair and enhancement. These findings underscore the importance of leveraging runtime insights to improve the robustness and accuracy of self-debugging methods.

## G Prompts

### G.1 Prompts for Initial Code Generation

The prompts for initial program generation are provided in Figure 8.

### G.2 Prompts for Test Generation

The prompts for test case generation are provided in Figure 9.

### G.3 Prompts for Debugging

The prompts for post-execution and in-execution self-debugging are provided in Figure 10 and Figure 11.

## Initial Generation (w/o self-debug)

```
===== SYSTEM =====  
You are an expert programming assistant.
```

```
===== USER =====  
Here is the given code to do completion:  
```python
```

```
def can_arrange(arr):  
    """Create a function which returns the largest index of an element which  
    is not greater than or equal to the element immediately preceding it. If  
    no such element exists then return -1. The given array will not contain  
    duplicate values.
```

Examples:

```
can_arrange([1,2,4,3,5]) = 3
```

```
can_arrange([1,2,3]) = -1
```

```
"""
```

```
```
```

Please complete the `can\_arrange` function and make sure that it matches the specification and passes all tests. You are not allowed to modify the given function signature. Think step by step and provide all completed codes in one code block.

```
===== ASSISTANT =====  
Sure, let's complete the `can_arrange` function step by step.
```

1. We need to iterate through the array from the second element to the last element.
2. For each element, we need to check if it is not greater than or equal to the element immediately preceding it.
3. If such an element is found, we return its index.
4. If no such element is found by the end of the loop, we return -1.

Here is the completed code:

```
```python  
def can_arrange(arr):  
    """Create a function which returns the largest index of an element which  
    is not greater than or equal to the element immediately preceding it. If  
    no such element exists then return -1. The given array will not contain  
    duplicate values.  
  
    Examples:  
    can_arrange([1,2,4,3,5]) = 3  
    can_arrange([1,2,3]) = -1  
    """  
    for i in range(len(arr) - 1, 0, -1):  
        if arr[i] < arr[i - 1]:  
            return i  
    return -1  
```
```

This function iterates through the array from the end to the beginning, checking if the current element is less than the previous element. If it finds such an element, it returns the index of that element. If no such element is found, it returns -1.

## In-Execution Self-Debugging w/ trace

```
===== USER =====
Given an input for the function `can_arrange([5, 4, 3, 2, 1])`, here is the code
execution trace block by block with the intermediate variable values:
[BLOCK-0]
# arr = [5, 4, 3, 2, 1]
for i in range(len(arr) - 1, 0, -1):
# arr = [5, 4, 3, 2, 1] i = 4
[BLOCK-1]
# arr = [5, 4, 3, 2, 1] i = 4
if arr[i] < arr[i - 1]:
# arr = [5, 4, 3, 2, 1] i = 4
[BLOCK-2]
# arr = [5, 4, 3, 2, 1] i = 4
return i
# arr = [5, 4, 3, 2, 1] i = 4 _ret = 4
```

Please explain the execution FOR EACH BLOCK and answer whether this block is correct or not based on the specifications and given samples in the problem.

```
===== ASSISTANT =====
Let's go through the execution trace block by block and explain the intermediate
variable values and the correctness of each block.
```

```
### [BLOCK-0]
```python
# arr = [5, 4, 3, 2, 1]
for i in range(len(arr) - 1, 0, -1):
# arr = [5, 4, 3, 2, 1] i = 4
```
**Explanation:**
- The array `arr` is `[5, 4, 3, 2, 1]`.
- The `for` loop is set to iterate from the last index (`len(arr) - 1`, which is 4) to the first index (1), decrementing by 1 each time.
- In the first iteration, `i` is initialized to 4.
**Correctness:**
- This block is correct. The loop is set up properly to iterate from the end of the array to the beginning.

### [BLOCK-1]
```python
# arr = [5, 4, 3, 2, 1] i = 4
if arr[i] < arr[i - 1]:
# arr = [5, 4, 3, 2, 1] i = 4
```
**Explanation:**
- The condition `if arr[i] < arr[i - 1]` is checked.
- For `i = 4`, `arr[4]` is 1 and `arr[3]` is 2.
- The condition `1 < 2` is true.
**Correctness:**
- This block is correct. The condition correctly checks if the current element is less than the previous element.
```

```

### [BLOCK-2]
```python
# arr = [5, 4, 3, 2, 1]   i = 4
return i
# arr = [5, 4, 3, 2, 1]   i = 4   _ret = 4
```

**Explanation:**
- Since the condition in [BLOCK-1] is true, the function returns the current index `i`, which is 4.

**Correctness:**
- This block is correct. The function correctly returns the index of the element that is not greater than or equal to the element immediately preceding it.

### Conclusion
The execution trace is correct and matches the expected behavior of the function based on the given specifications and samples. The function correctly identifies that the element at index 4 (value 1) is less than the element at index 3 (value 2) and returns the index 4.

```

### Post-Execution Self-Debugging w/ detail

```

===== USER =====
But it failed in the following test case:
assert can_arrange([5, 4, 3, 2, 1]) == 1, "Test Case 3: Descending order array"
# Real Execution Output: 4

Please fix the bug in the `can_arrange` function and make sure that the fixed code matches the specification and passes all tests. Provide the fixed code in one code block. If there is no bug in your previous code, repeat the code in a code block again. You are not allowed to modify the given function signature. Think step by step.

===== ASSISTANT =====
I see the issue. The function should return the largest index of an element which is not greater than or equal to the element immediately preceding it. In the case of a descending order array, the function should return the first index where this condition is met, which is `1` in the given test case.

To fix this, we need to iterate from the beginning of the array to the end, and return the first index where the condition is met. Here is the corrected code:

```python
def can_arrange(arr):
    """Create a function which returns the largest index of an element which is not greater than or equal to the element immediately preceding it. If no such element exists then return -1. The given array will not contain duplicate values.

    Examples:
    can_arrange([1,2,4,3,5]) = 3
    can_arrange([1,2,3]) = -1
    """
    for i in range(1, len(arr)):
        if arr[i] < arr[i - 1]:
            return i
    return -1
```

This function now correctly iterates from the beginning of the array to the end, and returns the first index where the current element is less than the previous element. This should pass all the given test cases.

```

Figure 7: An example with GPT-4o performing both post and in-execution self-debugging on a problem from HumanEval (HumanEval/135) respectively. Post-execution self-debugging wrongly corrects the program while in-execution self-debugging manages to keep the original answer.

### Code Generation Prompt for HumanEval

Here is the given code to do completion:

```
```python
{prompt}
```
```

Please complete the `{entry_point}` function and make sure that it matches the specification and passes all tests. You are not allowed to modify the given function signature. Think step by step and provide all completed codes in one code block.

### Code Generation Prompt for MBPP

Here is the given problem to solve:

```
```python
{prompt}
```
```

Please implement the `{entry_point}` function and make sure that it matches the specification and passes all tests. You are not allowed to modify the given function name and arguments in the test examples. Think step by step and provide all completed codes in one code block.

### Code Generation Prompt for Functional-Input Question in LiveCodeBench

Here is the given programming problem to solve:

```
{content}
```

Please generate a correct python program that matches the specification and passes all tests. Think step by step. You will use the following starter code to write the solution to the problem and enclose your code within delimiters.

```
```python
{starter_code}
```
```

### Code Generation Prompt for Stdin-Input Question in LiveCodeBench

Here is the given programming problem to solve:

```
{content}
```

Please generate a correct python program that matches the specification and passes all tests. Read the inputs from stdin solve the problem and write the answer to stdout (do not directly test on the sample inputs). Think step by step and enclose your code within delimiters as follows:

```
```python
# YOUR CODE HERE
```
```

Figure 8: Prompts for initial code generation (w/o *self-debug*).

### Test Generation Prompt for HumanEval

Here is the given code to do completion:

```
```python
{prompt}
```
```

Please provide ten comprehensive and valid test cases to verify whether the `{entry_point}` function correctly solves the problem. You are not allowed to implement the function. Think step by step and provide all test cases in one code block.

The format of test cases should be:

```
```python
assert {entry_point}(input) == expected_output, "Test Case Description"
```
```

### Test Generation Prompt for MBPP

Here is the given problem to solve:

```
```python
{prompt}
```
```

Please provide ten comprehensive and valid test cases to verify whether the `{entry_point}` function correctly solves the problem. You are not allowed to implement the function. Think step by step and provide all test cases in one code block.

The format of test cases should be:

```
```python
assert {entry_point}(input) == expected_output, "Test Case Description"
```
```

### Test Generation Prompt for LiveCodeBench

Here is the given programming problem to solve:

```
{content}
```

Please provide ten comprehensive test samples based on the specification and follow the format of the given sample.

Your response should be organized like below and no extra information is allowed (including explanation):

```
[Input]
<your input here>
[Output]
<your output here>
[Input]
...
```

Figure 9: Prompts for test generation.

### Post-Execution Debugging Prompt for HumanEval

```
{error}
```

Please fix the bug in the `{entry_point}` function and make sure that the fixed code matches the specification and passes all tests. Provide the fixed code in one code block. If there is no bug in your previous code, repeat the code in a code block again. You are not allowed to modify the given function signature. Think step by step.

### Post-Execution Debugging Prompt for MBPP

```
{error}
```

Please fix the bug in the `{entry_point}` function and make sure that the fixed code matches the specification and passes all tests. Provide the fixed code in one code block. If there is no bug in your previous code, repeat the code in a code block again. You are not allowed to modify the given function name and arguments in the test examples. Think step by step.

#### Post-Exec. Debugging Prompt for Functional-Input Questions in LiveCodeBench

{error}

Please fix the bug in the code and make sure that the fixed code matches the specification and passes all tests.  
You will use the following starter code to write the solution to the problem and enclose your code within delimiters.

```
```python
{starter_code}
```
```

#### Post-Exec. Debugging Prompt for Stdin-Input Questions in LiveCodeBench

{error}

Please fix the bug in the code and make sure that the fixed code matches the specification and passes all tests.  
Read the inputs from stdin solve the problem and write the answer to stdout (do not directly test on the sample inputs). Enclose your code within delimiters as follows.

```
```python
# YOUR CODE HERE
```
```

Figure 10: Prompts for post-execution debugging.

#### In-Execution Debugging Prompt

Given an input for the function `{test}`, here is the code execution trace block by block with the intermediate variable:  
{trace}

Please explain the execution FOR EACH BLOCK and answer whether this block is correct or not based on the specifications and given samples in the problem. If the program is correct, please restate it in one python code block. If it is incorrect, please fix the bug and provide the fixed code in a code block.

Figure 11: Prompt for in-execution self-debugging.