

# EvoGPT: Enhancing Test Suite Robustness via LLM-Based Generation and Genetic Optimization

Lior Broide<sup>a,\*,1</sup> and Roni Stern<sup>a,1</sup>

<sup>a</sup>Ben Gurion University of the Negev

## Abstract.

Large Language Models (LLMs) have recently emerged as promising tools for automated unit test generation. We introduce a hybrid framework called EvoGPT that combines LLM-based test generation with more traditional evolutionary search techniques to produce diverse, bug-revealing unit tests. Unit tests are initially generated using diverse temperature sampling to enhance behavioral and test suites diversity, followed by a generation-repair loop and coverage-guided assertion generation. The improved test suites are subsequently evolved using genetic algorithms, with a fitness function that prioritizes mutation score over traditional coverage metrics. This design reflects the primary goal of unit testing — detecting faults. Evaluated across multiple open-source Java projects, our framework achieves an average improvement of 10% in both code coverage and mutation score, compared to LLM and traditional search-based software testing baselines. These results demonstrate that combining LLM-driven diversity, targeted repair, and search-based optimization yields more effective and resilient test suites.

## 1 Introduction

Software testing is a fundamental practice in software engineering, critical for ensuring the reliability, correctness, and maintainability of software systems [6]. Among its various forms, *unit testing* serves as the first line of defense against bugs by verifying the behavior of individual program units in isolation. However, despite its importance, the process of manually designing high-quality unit tests remains labor-intensive and error-prone. Furthermore, the *oracle problem* i.e., the difficulty of defining the expected outcomes for arbitrary software behaviors, poses a persistent challenge, particularly when aiming to automate the generation of meaningful and verifiable test cases [21].

*Search-Based Software Testing* (SBST) is an extensively studied paradigm for automated test generation [13] [3]. Techniques such as evolutionary algorithms, implemented in tools like *EvoSuite* [9], systematically optimize test suites toward structural coverage objectives (e.g., branch and line coverage). While SBST methods often achieve high coverage, they typically generate test cases that lack readability, semantic relevance, insufficient fault detection, and their optimization processes are computationally expensive [20].

In recent years, *Large Language Models* (LLMs) have emerged as promising tools for automating unit test generation [7]. Approaches

such as *TestART* [12], *ChatUniTest* [8], and *A3Test* [2] demonstrate the ability of LLMs to produce syntactically correct, semantically coherent, and human-readable unit tests. These models leverage their natural language understanding and code generation capabilities to approximate human-level testing strategies. Nevertheless, LLM-based methods often suffer from limited exploration of complex input spaces, instability caused by hallucinations, and an inability to iteratively improve upon initial generations without external feedback [17].

Motivated by the complementary strengths and weaknesses of LLM-based and SBST-based approaches, we propose a *hybrid framework* that we call EvoGPT, which integrates diverse LLM-based generation with evolutionary search-based refinement. Our framework first generates a large, diverse pool of initial test suites by querying LLMs under varied temperature settings, encouraging exploration across diverse execution paths that include repair loops, and coverage guidance. Subsequently, an evolutionary algorithm refines the generated test suite through *selection*, *crossover* and *mutation* operations, optimizing for both coverage and mutation score. Empirical results show that our framework already significantly outperforms state-of-the-art LLM-only and SBST-only baselines in both code coverage and mutation score of successfully generated tests.

## Contributions

The main contributions of this paper are:

- **Hybrid Test Generation Framework:** a novel method for combining LLM-based test generation with evolutionary search to produce high-quality, high-coverage unit test suites.
- **Prompt- and Temperature-Driven Diversity:** We employ temperature-based sampling and different system prompt instructions during LLM querying to encourage diversity among initial test suites, improving input space exploration.
- **Effective evolutionary Optimization:** Our evolutionary loop achieves strong performance improvements, demonstrating the intrinsic effectiveness of hybridization.
- **Extensive Empirical Evaluation:** We rigorously evaluate our approach on standard benchmarks, demonstrating substantial improvements in code coverage and mutation testing over strong baselines, namely EvoSuite [9] and TestART [12].

## 2 Background and Problem Definition

A *unit test* is a self-contained function designed to validate the behavior of a specific method or unit of source code. As introduced in

\* Corresponding Author. Email: broidel@post.bgu.ac.il

<sup>1</sup> Equal contribution.

Test-Driven Development (TDD): "Tests are used to specify expected behavior and detect regressions when code changes over time" [6]. Executing a test typically involves invoking the target method under specific inputs and comparing the observed outcome against expected outputs through assertions. The result of a test is often binary: a test either *passes* or *fails*, depending on whether all assertions hold.

In the context of automated unit test generation, it is common to focus on testing a subset of methods within a class. Typically those that are externally visible and designed to be invoked directly. Following prior work [12], we designate the *public methods* of each class under test as **focal methods**. These focal methods serve as the primary targets for test generation and evaluation, allowing the testing effort to concentrate on actionable and externally observable behavior.

Evaluating the quality of a unit test involves several well-established metrics [9, 13, 15]:

- **Line Coverage of Correct Tests (LCCT)**: the proportion of source code lines in the focal methods that are executed at least once by a passing test. While high line coverage indicates broader code exploration, it does not guarantee semantic correctness.
- **Branch Coverage of Correct Tests (BCCT)**: the proportion of conditional branches (e.g., if-else, switch cases) in the focal methods that are exercised in passed tests. This metric captures control-flow depth and is more sensitive to logical path diversity.
- **Mutation Score of Correct Tests (MSCT)**: the percentage of artificially injected faults (mutants) in the focal methods that are detected, i.e., cause at least one test to fail. This is widely regarded as the strongest proxy for fault detection capability [9, 15].

The objective of an automated test generation algorithm is therefore to synthesize a suite of correct, diverse, and assertion rich tests that maximize test quality metrics, particularly LCCT, BCCT, and MSCT, while remaining readable and executable.

**Traditional test generation methods.** Several techniques have been proposed to automatically generate test suites. *Random test generation* methods such as Randoop [18] sample inputs and assertions without guided heuristics. While computationally inexpensive, they frequently produce shallow, brittle, or irrelevant tests that lack meaningful and relevant assertions, or practical utility. Search-Based Software Testing methods, such as EvoSuite [9] applies evolutionary algorithms to optimize the aforementioned test quality metrics over multiple generations. It begins by generating an initial population of random test suites, which are then iteratively refined using three core operators: *selection*, which favors test suites with higher fitness for reproduction; *crossover*, which combines parts of two parent test suites to create new ones; and *mutation*, which introduces small random changes to individual tests. This evolutionary process gradually improves coverage and fault detection ability over time. Applying evolutionary algorithms often requires implementing domain-specific fitness functions and cross-over and mutation operators. EvoSuite offers such an implementations for test generation. While EvoSuite and similar approaches can achieve high structural coverage, they struggle with assertion quality and test readability.

**LLM-based test generation methods.** More recently, *Large Language Models (LLMs)* have been used to synthesize test cases based on natural language prompts or code structure. Early approaches relied on single shot prompts (e.g., "Write a unit test for this method"), but more sophisticated techniques have since emerged. These include prompt engineering to guide assertion behavior, few-shot prompt-

ing to improve test structure, and iterative refinement using feedback such as stack traces or compilation errors.

A notable example is TestART [12], which combines LLM-based test generation with a feedback loop that repairs failing tests. In TestART, the model first generates a candidate test. If it fails to compile or run, the system captures the error. Then, it tries to fix it programmatically. If the error persists, it re-prompts the model with the error trace and original input, refining the test in multiple rounds. This repair loop significantly improves test pass rates (i.e., correct tests), and serves as an important foundation for our work. However, these LLM-generated tests often contain superficial validations and overfit to syntactic patterns rather than exercising the semantics of the code under test [23, 8, 19]. Moreover, they tend to neglect assertion diversity and fail to iteratively refine test quality in the absence of external feedback. Importantly, high code coverage does not necessarily imply high test effectiveness. Empirical studies have shown that tests with high line or branch coverage may still fail to detect subtle faults if they lack strong or diverse assertions. This further motivates the use of mutation score as a more reliable adequacy criterion [14]. This paper focuses on leveraging LLMs in conjunction with evolutionary algorithms to achieve this goal through multi-agent diversity, iterative test repair, and evolutionary optimization.

### 3 The EvoGPT System

In this section, we describe EvoGPT, a hybrid framework that combines Large Language Models (LLMs) and evolutionary search to generate and optimize robust test suites. Figure 1 provides a high-level illustration of EvoGPT. EvoGPT applies a two-phase process: initial diversification via prompt-engineered LLM agents, followed by evolutionary refinement. Prior to test generation, it perform lightweight preprocessing on the source code to reduce token complexity and context window. Specifically, we remove excessive or lengthy inline comments, documentation blocks, and unreachable code segments that could overload the model's context window and negatively affect generation quality [10]. Then, an LLM-based component is used to generate an initial pool of diverse test suites.

Specifically, we instantiate five distinct LLM agents, each defined by a different prompting strategy (e.g., focusing on edge cases, deep object chains, or creative assertion styles). Each agent is asynchronously queried five times, resulting in a total of 25 diverse and structurally unique test suites. During this generation process, a lightweight repair loop is applied to each test suite: when a generated test fails due to compilation or runtime errors, the system captures the stack trace and re-invokes the LLM with this context to produce a fixed version. This helps improve test validity and robustness early in the pipeline.

The resulting 25 test suites are passed to the Genetic Algorithm (GA) component, serving as its initial population. The GA component begins a search-based optimization process in which it iteratively evolves this population by applying selection, crossover and mutation operations guided by fitness function. Next, we describe in details the main components of EvoGPT: the LLM component and the GA component.

#### 3.1 LLM Component

To generate tests for a given class, we instantiate 25 LLM agents, each 5 LLM agents configured with a different temperature setting and a corresponding system prompt as seen in Table 1. The temperature value is used to control the stochasticity of the language

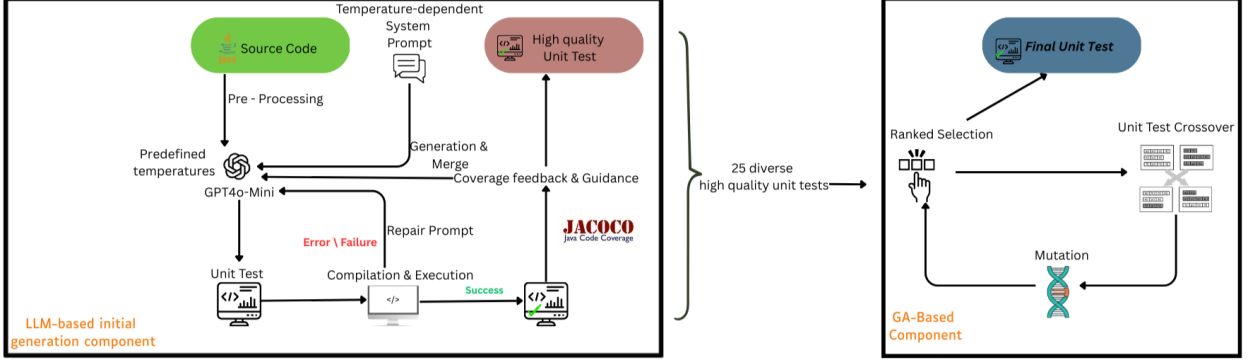


Figure 1. Illustration of the EvoGPT system.

model output. Prior research has shown that varying temperature values leads to a broader range of generated outputs, exposing different behavioral patterns and increasing semantic variance across samples [1]. To complement this, each agent is given a custom system prompt, which directs it to approach the test suite generation task from a slightly different perspective (e.g., focusing on execution paths, edge values, or assertion robustness). Table 1 lists for each of the five types of LLM agents we used their temperature value, a qualitative description of their prompt (column “Prompt Strategy”), and intended purpose (column “Purpose”). The exact system prompts used for each LLM agent are included in the supplementary material. This diversity promotes both semantic variability and robustness in the generated tests, and was refined empirically during early experimentation.

Agent	Temperature	Prompt Strategy	Purpose
A1	0.3	Standard unit testing	Reliable structure
A2	0.6	Emphasize assertion diversity	Broader observation
A3	0.8	“Try hard” creative agent	Path and value exploration
A4	0.5	Focus on edge conditions	Boundary case detection
A5	0.4	Uses long object chains	Deeper semantic chains

Table 1. LLM agent configuration: temperature and prompt strategy paired to encourage diversity and robustness. Full prompts are provided in the supplementary material.

**Example 1.** To illustrate the effect of prompt-based diversity, we analyzed test suites generated by our five agents (A1-A5) for the focal method `toJsonTree()` in the `Gson` library. Specifically, we sampled test generated with agents A2, A3, and A5 (“standard unit testing”, “emphasize assertion diversity”, and “use long object chains”). Although all agents target the same method, each produces semantically distinct test logic. Table 2 provides a qualitative summary of their key behaviors. For example, Agent A2 emphasizes assertion richness, including non-nullity and comparison against `JsonNull.INSTANCE`, Agent A4 focuses on string based edge inputs and value equality, and Agent A5 constructs a structured map input and checks the JSON object’s type. Full versions of the generated tests are included in the supplementary material.

These variations reflect the agents’ distinct prompting strategies, supporting our design choice to introduce semantic diversity through prompt engineering. While our agents also differ in temperature, we primarily attribute the observed diversity to the prompt styles. Preliminary experiments suggested that combining temperature varia-

tion with prompt diversity may further increase behavioral variance across tests.

Table 2. Qualitative summary of diverse test behaviors for `toJsonTree()` by agent type

Agent ID	Key Code / Behavior
A2	Includes assertions for non-nullity and inequality against <code>JsonNull.INSTANCE</code>
A4	Inputs string literal; asserts equality with "test" using <code>getAsString()</code>
A5	Uses <code>Map&lt;String, String&gt;</code> ; asserts object via <code>isObject()</code>

After generating the initial set of tests, each LLM agent performs a generation-repair loop where it iteratively refines its test case in response to runtime feedback such as compilation errors or failing assertions. Such a generation-repair process was shown to enhance test validity and robustness in prior LLM-based frameworks such as ChatUnitTest [8], TestART [12] and TestPilot [19]. In our system, this loop leverages a simplified version of TestART’s repair strategy: when a generated test fails with a runtime error or stack trace, we programmatically extract the stack trace and invoke a minimal test patching heuristic to remove or correct the failing statement. This step addresses the tendency of LLM-generated tests to fail due to repetitive, structurally similar errors, such as unresolved imports or assertion mismatches.

Rather than relying on additional LLM calls—which often perpetuate these failures, our minimal repair heuristic directly fixes common runtime and compilation errors, following TestART’s insight that targeted, template-based repairs are more stable and cost-effective for improving test pass rates without iterative hallucination-prone feedback loops [12]. In case these programmatic fixes do not work, the LLM is injected with a repair prompt, along with the stack trace, defining the runtime or compilation error. The exact repair prompt is listed in the supplementary material. This repair loop is repeated 4 times to prevent the context window from being overloaded [10]. If errors in the tests persist, the test methods responsible for these errors are automatically removed from the test suite.

Once a valid test suite is synthesized, it is compiled and executed using JaCoCo, a Java Code Coverage tool, which provides detailed line and branch coverage information. The resulting coverage report is parsed and passed to the next module: the Coverage Analysis

Agent. This agent uses the same temperature as the original generation agent to maintain behavioral consistency. It is then instructed to generate additional complementary tests, aimed specifically at covering the missed branches and lines reported in the coverage report. This strategy was found valuable for enhancing code coverage, as demonstrated in TestART [12]. Finally, the tests produced by both the initial generation agent and its corresponding coverage enhancement process are merged into a single composite test suite. This forms one member of the initial population in our evolutionary algorithm.

## 3.2 Evolutionary Component

---

### Algorithm 1 EvoGPT Evolutionary Optimization Loop

---

```

1: Input: Initial LLM-generated test population
2: Output: Highest-fitness test suite
3: Compute initial fitness for all individuals
4: Store best initial individual
5: while time budget not exceeded do
6:   Initialize empty elite pool
7:   while elite pool not full do
8:     if elite pool has perfect test suite then
9:       Return best from elite pool
10:    end if
11:    Select two parents via ranked selection
12:    if random() <  $p_c$  then
13:      Perform crossover → two offspring
14:    else
15:      Copy parents as offspring
16:    end if
17:    Apply assertion mutation to each offspring
18:    Compute fitness for both offspring
19:    if offspring better or equal (and shorter) than parents then
20:      Add offspring to elite pool
21:    else
22:      Add parents to elite pool
23:    end if
24:  end while
25:  Replace population with elite pool
26: end while
27: if best initial not in population then
28:   Add it to population
29: end if
30: return highest-fitness test suite

```

---

To refine and evolve the test suite population generated by the LLM agents, we employ the tailored genetic algorithm (GA) listed in Algorithm 1. Our GA component builds on these prior works, using similar crossover, mutation, and selection operators for Java test suite files, coupled with a fitness evaluation grounded in code coverage and mutation resistance. The GA operates over a fixed time budget and evolves a population of test suites toward higher robustness and quality. The algorithm terminates early if a chromosome (i.e., test suite, in EvoGPT’s context) reaches a perfect fitness score (which is 100 in the defined below fitness function that we used). If the initial population pool has no perfect test suite, the GA begins, and returns the highest scoring candidate at the end of evolution.

### 3.2.1 Fitness Score

The fitness function plays a central role in guiding the GA toward higher-quality solutions through what is known as *evolutionary pressure* - the selective process by which better-performing candidates are favored for reproduction, gradually improving the population over successive generations. In EvoGPT, the primary goal is to generate tests that are effective at exposing faults. To that end, our fitness function includes the *mutation score of correct tests* (MSCT), which measures a test’s ability to detect artificially injected faults and is widely regarded as a strong proxy for fault-detection capability [9]. To ensure that structural exploration is not neglected, we also incorporate branch coverage (BCCT) and line coverage (LCCT) into the fitness function. The final fitness score is defined as a weighted combination of these three components:

$$\text{fitness} = 0.3 \cdot \text{BCCT} + 0.2 \cdot \text{LCCT} + 0.5 \cdot \text{MSCT} \quad (1)$$

This composite formulation encourages the generation of test suites that not only explore diverse execution paths but also robustly detect behavioral anomalies.

**Coefficient Rationale.** The specific weights were chosen based on our guiding priorities: MSCT receives the highest weight (0.5) to emphasize fault detection; BCCT (0.3) is prioritized over LCCT (0.2) as it better captures logical path diversity. While these coefficients were informed by prior literature [9] and early empirical experiments, they may be further tuned and experimented in future work to suit specific domains or testing objectives.

### 3.2.2 Selection

In our evolutionary framework, we employ ranked selection to choose parent test suites for reproduction. This method involves sorting the current population based on fitness scores and assigning selection probabilities according to these ranks rather than raw fitness values. Ranked selection is advantageous [22] as it maintains genetic diversity within the population and mitigates the risk of premature convergence [5], a scenario where the algorithm converges to suboptimal solutions early in the search process. By preventing highly fit individuals from dominating the gene pool too quickly, ranked selection ensures a more balanced exploration of the solution space [11].

This approach provides a controlled *selection pressure*. High selection pressure can accelerate convergence but may lead to a loss of diversity, while low selection pressure promotes exploration but may slow down progress. Ranked selection strikes a balance by giving moderately fit individuals a non-negligible chance to reproduce, thus allowing less fit test suites to contribute to the next generation. This balance between exploration and exploitation is crucial for the robustness and effectiveness of genetic algorithms in complex optimization tasks.

### 3.2.3 Crossover

In genetic algorithms, one of the primary operators that drive population evolution is **crossover**. Crossover combines two parent individuals to produce two offsprings, allowing promising traits to propagate and mix across generations.

In our implementation, we use a *crossover probability* of 0.8. That is, with an 80% chance, the algorithm selects two parents and performs a structural crossover to generate two offsprings.

This crossover operator blends the internal structure of the two parent test suites—specifically their test methods, helper methods, and inner classes. One offspring inherits approximately 80% of its test methods from Parent 1 and 20% from Parent 2, while the second offspring receives the inverse proportion.

### 3.2.4 Mutation

An additional operator used in genetic algorithms is **Mutation**. Mutation introduces random variations into individuals, ensuring continued diversity and helping escape local optima. We introduce a novel mutation strategy inspired by recent findings in LLM-guided test enhancement [4]. Instead of traditional statement insertion, deletion, or alteration [9], we augment each offspring’s test methods with additional assertions, generated by an LLM agent. The assertion generating agent operated with a relatively high temperature of 0.5, to generate diverse and creative assertions. The exact prompt is listed in Figure 3.2.4.

These assertions aim to improve the meaningfulness and relevance of test validations by going beyond shallow checks, e.g., “!=null”, to more context aware conditions that verify object properties or behavior. In doing so, they increase the likelihood that the test will detect subtle faults introduced by code mutations, thereby improving the test’s fault detection capability and robustness.

Each test method has a mutation probability defined as:

$$P_{\text{mutation}} = \frac{1}{N_{\text{tests}}} \quad (2)$$

Where  $N_{\text{tests}}$  is the number of test methods in the mutated test suite class. Thus, a single mutation is performed in every test suite class in expectation. This conservative rate is inspired by common evolutionary computation practice, and aims to preserve test structure, avoid overspecialization to current code conditions, and encourage controlled semantic diversity.

## 4 Experimental Results

We performed an extensive experimental evaluation to assess the performance of EvoGPT. We used **Defects4J** [16] as our primary benchmark, a well-established repository of real-world Java projects with reproducible bugs and accompanying test suites. It has been widely used in prior work on automated testing and program repair due to its diversity, complexity, and grounding in real software engineering practice. For our experiments, we select a representative subset of classes from multiple projects (e.g., Gson, Lang, CLI, CSV). We extracted public and non-abstract classes from these four projects as focal classes and extracted their public methods as focal methods.

**Table 3.** Dataset overview: Projects and their focal method counts

Dataset	Project name	Abbr.	Version	Focal methods
Defects4J	Gson	Gson	2.10.1	378
	Commons-Lang	Lang	3.1.0	1728
	Commons-Cli	Cli	1.6.0	177
	Commons-Csv	Csv	1.10.0	137

To measure the quality of the generated tests, we used the three test quality metrics mentioned above – LCCT, BCCT, and MSCT – restricted to passing tests and focal method coverage. Mutation metrics are based on the mutant analysis of PITest (version 1.19.0), while coverage is computed with JaCoCo (version 0.8.12). A time budget

```
You are a mutation agent for evolutionary unit test generation.

Your role is to enhance the robustness of Java unit test methods by intelligently adding new assertions. These additional assertions should verify more properties of the system under test (SUT), including outputs, internal state (via getters), object properties, or observable side effects.

Follow these strict guidelines:

1. Preserve the original logic of the test. Do not modify or remove existing code.
2. Add 1 to 5 new assertions that enhance the strength of the test by:
   - Checking additional return values.
   - Inspecting object state using getters or fields.
   - Validating side effects (e.g., object creation, mutation, collections modified).
3. Prefer using existing variables and objects in the test.
4. If a needed object is missing, you may instantiate new objects within the test method, but only if:
   - They are directly relevant to additional assertions.
   - They do not significantly change the original test’s intent.
   - The object is imported into the unit test’s java file.
5. You may use helper methods (e.g., 'assertEquals', 'assertTrue', 'assertNotNull', etc.).
6. Maintain consistent code style and indentation.
7. Do not explain the code. Only output the updated Java test method.

You will be provided:

- The original Java test method to mutate.
- The source code of the class under test.
- Optional helper methods or imported utility code used in the test.

Your output must be the full, modified Java test method, with @Test decorator, with the added assertions inserted logically and consistently within the method body. Output the test method code without markdown formatting (``java`).
```

**Figure 2.** Prompt for generating a mutation.

of 300 seconds was allocated for every iteration of the GA, and the LLM we used in our experiments is gpt-4o-mini.

### 4.1 Comparison Against Baselines

Project	Framework	LCCT	BCCT	MSCT	# of Tests
Gson	EvoGPT	<b>94.1</b>	<b>90.6</b>	<b>89.1</b>	<b>284</b>
	TestART	81.9	79.4	76.3	292
	EvoSuite	81.1	79.2	67.1	803
Lang	EvoGPT	<b>95.6</b>	<b>94.2</b>	<b>87.2</b>	848
	TestART	87.7	84.4	82.7	<b>825</b>
	EvoSuite	90.1	86.8	75.2	2729
CSV	EvoGPT	<b>96.2</b>	<b>94.1</b>	<b>98.4</b>	133
	TestART	80.7	79.4	76.7	<b>116</b>
	EvoSuite	81.7	70.0	62.0	546
CLI	EvoGPT	<b>96.0</b>	<b>95.4</b>	<b>91.0</b>	<b>161</b>
	TestART	87.0	86.4	85.8	182
	EvoSuite	93.8	93.9	85.8	535
Total	EvoGPT	<b>95.5</b>	<b>93.6</b>	<b>91.4</b>	1426
	TestART	84.3	82.4	80.4	<b>1415</b>
	EvoSuite	86.7	82.5	72.5	4613

**Table 4.** Framework comparison across benchmarks

Table 4 presents the average LCCT, BCCT, and MSCT scores

across our four benchmark projects (Gson, Lang, CSV, CLI). The right-most column shows the summed number of test generated by each method. EvoGPT is compared to two strong baselines: EvoSuite [9], a widely used search-based test generation tool, and TestART [12], a recent state of the art framework for LLM-based unit test generation. EvoSuite represents the traditional SBST approach for test generation and TestART represents more modern LLM-driven techniques. We chose TestART since it has been shown to outperform previous LLM-based approaches such as ChatUniTest [8].

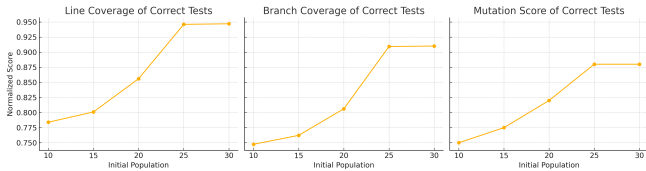
As can be seen in Table 4, EvoGPT consistently outperforms both baselines across all metrics. It achieves the highest average mutation score (0.91), improving over TestART (0.86) and EvoSuite (0.80). Similarly, EvoGPT leads in both LCCT and BCCT, suggesting that its multi-agent diversity and assertion-enhancement mutation improve both path exploration and semantic validation. These improvements are especially notable in challenging domains like CSV and CLI, where LLM-generated tests alone often underperform.

## 4.2 Sensitivity Analysis and Ablation Studies

To better understand the contribution of individual components in EvoGPT, we conduct a series of parameter sensitivity analysis and ablation studies. These experiments systematically disable or modify core aspects of the framework to assess their impact on test effectiveness. In particular, we evaluate the effect of (1) the initial population size, (2) the presence of our assertion-enhancing mutation strategy, and (3) the diversity introduced by varying temperature values across LLM agents. Each setting is evaluated using the same metrics as in prior sections (LCCT, BCCT, MSCT), averaged across all benchmarks.

### 4.2.1 Impact of Initial Population Size

To evaluate the role of LLM-driven diversity, we varied the number of initial agents from 10 to 30, while keeping all other parameters fixed. Figure 3 shows the resulting trends in LCCT, BCCT, and MSCT.



**Figure 3.** Impact of initial LLM-generated test population size on LCCT, BCCT, and MSCT.

We observe a clear monotonic gain up to a population size of 25, after which scores plateau.

### 4.2.2 Ablation Studies

TO quantify the impact of the core component in EvoGPT, we experimented with different versions of EvoGPT in which some of its components are disabled. The results are listed in Table 5. As can be seen, disabling the evolutionary optimization (i.e., the GA component) entirely leads to the steepest performance degradation, with LCCT, BCCT, and MSCT dropping to 84.2%, 83.4%, and 81.1% respectively. This underscores the central role of the GA component

**Table 5.** Ablation study: Effect of mutation and temperature diversity on EvoGPT

Configuration	LCCT (%)	BCCT (%)	MSCT (%)
EvoGPT w/o evolutionary opt.	84.2	83.4	81.1
EvoGPT w/o temp. diversity	87.4	86.0	83.0
EvoGPT w/o Mutation	94.1	92.5	89.4
<b>EvoGPT (Full)</b>	<b>95.5</b>	<b>93.6</b>	<b>91.4</b>

in refining the test suite beyond the initial LLM-generated population. Removing temperature diversity and using only standard unit testing prompt, results in noticeable declines as well, most significantly in MSCT (down to 83%), which reflects a narrower exploration of assertion and input behaviors when generation is less varied. Disabling the mutation operator in the GA component yields a smaller but still measurable drop in effectiveness, particularly in mutation score (from 91.4% to 89.4%), indicating that assertion enrichment contributes to fault detection strength. Collectively, these results demonstrate that each component—evolutionary refinement, generation diversity, and assertion mutation—adds value, and that the full EvoGPT configuration achieves the most effective test outcomes.

## 4.3 Limitations

While EvoGPT generates decent results across the chosen benchmarks, our approach has some limitations.

- **Runtime.** One notable limitation of our framework lies in its computational cost. Although the initial population generation is executed asynchronously, generating 25 diverse unit test candidates using OpenAI’s API remains time-consuming, particularly for large and complex classes. This step introduces non-trivial latency due to API response times. Furthermore, the evolutionary optimization phase, despite being efficient in navigating the search space, also incurs a significant runtime per generation, as each candidate must be evaluated through instrumentation, compilation, and coverage and mutation analysis. As a result, the end-to-end test generation process for a single source class can be relatively slow, which may limit scalability in large-scale or time-constrained testing scenarios.
- **Monetary cost.** Our framework relies heavily on LLM-based generation throughout the test synthesis process: initializing with 25 unique test suites per class and further invoking the model for assertion enhancement during mutation. While this design maximizes test diversity and semantic richness, it comes with a tangible monetary cost due to many API calls to commercial LLM providers (OpenAI in our case). As such, the per-class cost of test generation may be prohibitive at scale, especially in industrial settings with hundreds or thousands of classes.

## 4.4 Threats to validity

**Internal Validity.** Our reported coverage and mutation scores depend on the accuracy of third-party tools: JaCoCo for line/branch coverage, and PITest for mutation analysis. These tools are widely used, but they are not perfect. For example, PITest may not detect equivalent mutants (mutants that behave identically to the original code), which can falsely lower the mutation score. Similarly, JaCoCo may miss certain paths or lines due to bytecode instrumentation quirks. These limitations could slightly distort the results and should be considered when interpreting our metrics.

**Construct Validity.** We evaluated generated test suites fitness using structural and mutation-based metrics (LCCT, BCCT, MSCT), but do not measure other qualitative aspects such as readability, assertion relevance, or developer trust. Moreover, our mutation strategy improves assertion robustness, but it may not always align with human-written test intent. In addition, EvoGPT leverages temperature-based LLM sampling, which are stochastic. Thus, the generated test suites may vary across runs. While our reported results reflect representative averages, future work could analyze variance across seeds to better characterize the stability of the generated test suites.

**External Validity.** Our experiments are limited to four projects from Defects4J. Although these span a range of domains, they may not fully represent the complexity of industrial codebases. Furthermore, we focused only on public methods (our focal methods), which may not capture the full interaction behavior in production scenarios.

**Reproducibility.** As noted above, EvoGPT relies on LLM outputs which are stochastic. Moreover, exact reproducibility depends on access to specific LLMs (e.g., gpt-4o-mini) and API configurations. We aim to release all prompts, scripts, and evaluation data to support replication.

## 5 Conclusion and Future Work

In this work, we introduced **EvoGPT**, a hybrid test generation framework that combines temperature-diverse LLM agents with evolutionary optimization and a novel mutation strategy centered on assertion enhancement. Our system demonstrates consistent improvements over established baselines in terms of line coverage, branch coverage, and mutation score across several real-world Java projects.

By leveraging multi-agent diversity and evolutionary refinement, EvoGPT effectively bridges the gap between high-level language model capabilities and rigorous structural test optimization. EvoGPT’s integration of LLM-generated diversity with evolutionary refinement positions it as a practical tool for enhancing automated test generation in real-world software development environments. Its ability to produce robust and semantically rich test cases can help developers identify defects more efficiently, thereby improving software quality and reliability.

Future work will expand the experimental scope to include a broader selection of projects from the Defects4J benchmark and benchmarks in other programming languages. In addition, we aim to assess the *bug-detection capability* of EvoGPT-generated test suites by evaluating whether they expose real historical defects. This can be measured by running EvoGPT on a version committed after some bug is fixed and then running the generated test suite using the faulty version, checking if the generated assertions “capture” the bug.

In addition, we plan to test the scalability and robustness of EvoGPT on large-scale, production-grade open-source repositories, where challenges such as complex class hierarchies, third-party dependencies, and flaky behavior are more prevalent. These extensions will allow us to further validate the practical utility of EvoGPT in software testing workflows and to better understand the limitations of current LLM-driven testing approaches in realistic settings. Another direction for future work is to integrate EvoGPT into a lifelong testing framework, which includes choosing intelligently which tests should be run and for which classes should we generate more test after every revision.

## References

- [1] A. Agarwal, K. Mittal, A. Doyle, P. Sridhar, Z. Wan, J. A. Doughty, J. Savelka, and M. Sakr. Understanding the role of temperature in diverse question generation by gpt-4. In *Proceedings of the 55th ACM Technical Symposium on Computer Science Education V. 2, SIGCSE 2024*, page 1550–1551, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400704246. doi: 10.1145/3626253.3635608. URL <https://doi.org/10.1145/3626253.3635608>.
- [2] S. Alagarsamy, C. Tantithamthavorn, and A. Aleti. A3test: Assertion-augmented automated test case generation. *Information and Software Technology*, 176:107565, 2024.
- [3] H. Almulla and G. Gay. Learning how to search: generating effective test cases through adaptive fitness function selection. *Empirical Software Engineering*, 27(2):38, 2022.
- [4] N. Alshahwan, J. Chheda, A. Finogenova, B. Gokkaya, M. Harman, I. Harper, A. Marginean, S. Sengupta, and E. Wang. Automated unit test improvement using large language models at meta. In *Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering*, pages 185–196, 2024.
- [5] J. E. Baker. Adaptive selection methods for genetic algorithms. In *Proceedings of the first international conference on genetic algorithms and their applications*, pages 101–106. Psychology Press, 2014.
- [6] K. Beck. *Test driven development: By example*. Addison-Wesley Professional, 2022.
- [7] M. Chen, J. Twarek, H. Jun, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [8] Y. Chen, Z. Hu, C. Zhi, J. Han, S. Deng, and J. Yin. Chatunitest: A framework for llm-based test generation. In *Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering*, pages 572–576, 2024.
- [9] G. Fraser and A. Arcuri. Whole test suite generation. *IEEE Transactions on Software Engineering*, 39(2):276–291, 2013.
- [10] M. Gao, T. Lu, K. Yu, A. Byerly, and D. Khashabi. Insights into LLM long-context failures: When transformers know but don’t tell. In Y. Al-Onaizan, M. Bansal, and Y.-N. Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7611–7625, Miami, Florida, USA, Nov. 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.447. URL <https://aclanthology.org/2024.findings-emnlp.447/>.
- [11] D. E. Goldberg and K. Deb. A comparative analysis of selection schemes used in genetic algorithms. In *Foundations of genetic algorithms*, volume 1, pages 69–93. Elsevier, 1991.
- [12] S. Gu, Q. Zhang, C. Fang, F. Tian, L. Zhu, J. Zhou, and Z. Chen. Testart: Improving llm-based unit testing via co-evolution of automated generation and repair iteration. *arXiv preprint arXiv:2408.03095*, 2024.
- [13] M. Harman. The current state and future of search based software engineering. In *Future of Software Engineering (FOSE’07)*, pages 342–357. IEEE, 2007.
- [14] L. Inozemtseva and R. Holmes. Coverage is not strongly correlated with test suite effectiveness. In *Proceedings of the 36th international conference on software engineering*, pages 435–445, 2014.
- [15] Y. Jia and M. Harman. An analysis and survey of the development of mutation testing. *IEEE transactions on software engineering*, 37(5): 649–678, 2010.
- [16] R. Just, D. Jalali, and M. D. Ernst. Defects4j: A database of existing faults to enable controlled testing studies for java programs. In *Proceedings of the 2014 international symposium on software testing and analysis*, pages 437–440, 2014.
- [17] A. Madaan, N. Tandon, P. Gupta, S. Hallinan, L. Gao, S. Wiegreffe, U. Alon, N. Dziri, S. Prabhunoye, Y. Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594, 2023.
- [18] C. Pacheco and M. D. Ernst. Randoop: feedback-directed random testing for java. In *Companion to the 22nd ACM SIGPLAN conference on Object-oriented programming systems and applications companion*, pages 815–816, 2007.
- [19] M. Schäfer, S. Nadi, A. Eghbali, and F. Tip. An empirical evaluation of using large language models for automated unit test generation. *IEEE Transactions on Software Engineering*, 50(1):85–105, 2024. doi: 10.1109/TSE.2023.3334955.
- [20] N. Setiani, R. Ferdiana, and R. Hartanto. Test case understandability model. *IEEE Access*, 8:169036–169046, 2020.
- [21] E. J. Weyuker. On testing non-testable programs. *The Computer Journal*, 25(4):465–470, 1982.
- [22] D. Whitley. The genitor algorithm and selection pressure: Why rank-based allocation of reproductive trials is best. In *Proceedings of the third international conference on Genetic algorithms*, 1989.

- [23] Z. Yuan, M. Liu, S. Ding, K. Wang, Y. Chen, X. Peng, and Y. Lou. Evaluating and improving chatgpt for unit test generation. *Proceedings of the ACM on Software Engineering*, 1(FSE):1703–1726, 2024.