# KTester: Leveraging Domain and Testing Knowledge for More Effective LLM-based Test Generation

Anji Li
lianj8@mail2.sysu.edu.cn
Sun Yat-sen University
Zhuhai, China

Mingwei Liu
liumw26@mail.sysu.edu.cn
Sun Yat-sen University
Zhuhai, China

Zhenxi Chen
chenzhx236@mail2.sysu.edu.cn
Sun Yat-sen University
Zhuhai, China

Zheng Pei
peizh3@mail2.sysu.edu.cn
Sun Yat-sen University
Zhuhai, China

Zike Li
lizk8@mail2.sysu.edu.cn
Sun Yat-sen University
Zhuhai, China

Dekun Dai
daidk@mail2.sysu.edu.cn
Sun Yat-sen University
Zhuhai, China

Yanlin Wang
wangylin36@mail.sysu.edu.cn
Sun Yat-sen University
Zhuhai, China

Zibin Zheng
zhzibin@mail.sysu.edu.cn
Sun Yat-sen University
Zhuhai, China

## Abstract

Automated unit test generation using large language models (LLMs) holds great promise but often struggles with generating tests that are both correct and maintainable in real-world projects. This paper presents KTester, a novel framework that integrates project-specific knowledge and testing domain knowledge to enhance LLM-based test generation. Our approach first extracts project structure and usage knowledge through static analysis, which provides rich context for the model. It then employs a testing-domain-knowledge-guided separation of test case design and test method generation, combined with a multi-perspective prompting strategy that guides the LLM to consider diverse testing heuristics. The generated tests follow structured templates, improving clarity and maintainability. We evaluate KTester on multiple open-source projects, comparing it against state-of-the-art LLM-based baselines using automatic correctness and coverage metrics, as well as a human study assessing readability and maintainability. Results demonstrate that KTester significantly outperforms existing methods across six key metrics, improving execution pass rate by 5.69% and line coverage by 8.83% over the strongest baseline, while requiring less time and generating fewer test cases. Human evaluators also rate the tests produced by KTester significantly higher in terms of correctness, readability, and maintainability, confirming the practical advantages of our knowledge-driven framework.

## 1 Introduction

Unit testing is a fundamental practice in software development that verifies whether a method behaves as expected, playing a key role in ensuring software correctness, maintainability, and enabling regression testing [1]. As the first line of defense in the development lifecycle, it helps detect and localize bugs early, preventing their propagation [2]. However, writing high-quality unit tests manually is often time-consuming and error-prone, as it requires developers to deeply understand the method's behavior, construct valid input states, and define precise assertions [3].

For a method under test (*i.e.,* often called as the focal method), its well-constructed unit test typically consists of two parts: the test prefix, which sets up the necessary objects and environment to bring the system into a testable state, and the test oracle, which verifies that the actual output matches the expected behavior [4]. High-quality unit tests not only improve software reliability but also serve as living documentation, enabling easier code comprehension and maintenance [5].

To reduce manual effort, automated test generation techniques have been proposed to generate a suite of unit tests with the main goal of maximizing the coverage in the software under test. Traditional unit test generation techniques include search-based [6–8], constraint-based methods [9–11], and random-based strategies [12, 13]. While effective in achieving code coverage, the generated tests are often hard to read and maintain, limiting practical adoption [14].

Recently, large language models (LLMs) have shown great promise in generating more human-like test code. These LLM-based methods demonstrate promising capabilities in understanding code semantics and synthesizing test code without explicit test specifications. Despite this potential, current LLM-based methods still suffer from key limitations that hinder their practical adoption due to a fundamental lack of essential knowledge. In particular, these models often lack access to **project-specific knowledge**, such as how to correctly instantiate and use classes, or how utility methods and APIs interact across modules. They also overlook **testing domain knowledge**, including core principles like boundary value analysis, exception handling, or the separation between test design and implementation. **This lack of both project-specific and testing-domain knowledge often results in test code that is unreliable and difficult to maintain.**

As illustrated in Figure 1 to 3, these issues commonly arise due to the absence of project-specific knowledge and test domain expertise in existing LLM-based test generation methods. For instance, Figure 1 shows a test that incorrectly constructs the *labels* Map used by *SparkApplication*, inserting the key "dragName" instead of the correct constant *Constants.DAG_NAME_LABE*. This error arises from the test's lack of awareness of constructor dependencies

and configuration requirements. Figure 2 depicts a test with insufficient assert statements, weakening its effectiveness. Meanwhile, Figure 3 shows tests that are inconsistently structured, relying on hard-coded values and missing essential setup logic, which results in fragile and difficult-to-maintain code.

To overcome limitations in existing LLM-based unit test generation, we propose a novel framework that integrates project-specific knowledge with software testing domain knowledge to guide the test generation process. Our approach consists of two main phases: an offline knowledge extraction phase and an online test generation pipeline. In the offline phase, we perform static analysis on the target project to extract two key types of knowledge: project structure knowledge, which includes class definitions, method signatures, and field declarations; and project usage knowledge, which covers invocation patterns, dependencies, and related functions. This comprehensive project knowledge is used to provide rich contextual information that improves the relevance and accuracy of the generated tests. The online test generation pipeline involves five sequential steps: (1) test class framework generation, which sets up reusable scaffolding such as setup and teardown methods; (2) multi-perspective test case design, where test scenarios are created from complementary testing heuristics; (3) test method transformation, which converts structured test cases into executable test methods; (4) test class integration, assembling all generated methods into a coherent test class; and (5) test class refinement, which improves code quality and maintainability. Throughout these steps, the pipeline leverages both the extracted project knowledge and guidance from established software testing principles.
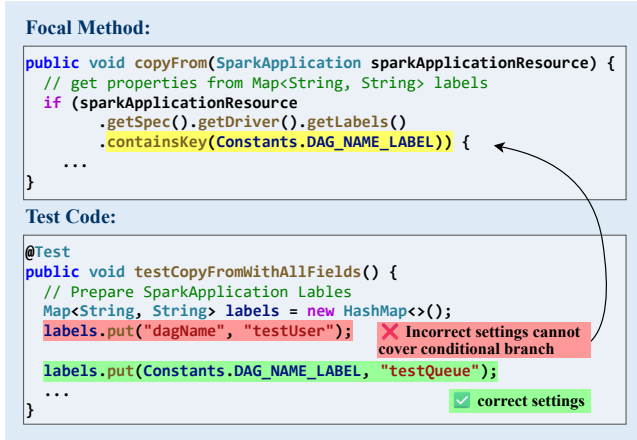
**Focal Method:**
```java
public void copyFrom(SparkApplication sparkApplicationResource) {
  // get properties from Map<String, String> labels
  if (sparkApplicationResource
      .getSpec().getDriver().getLabels()
      .containsKey(Constants.DAG_NAME_LABEL)) {
    ...
  }
}
```

**Test Code:**
```java
@Test
public void testCopyFromWithAllFields() {
  // Prepare SparkApplication Lables
  Map<String, String> labels = new HashMap<>();
  labels.put("dagName", "testUser");   ❌ Incorrect settings cannot
                                          cover conditional branch
  labels.put(Constants.DAG_NAME_LABEL, "testQueue");
  ...                                  ✅ correct settings
}
```

**Figure 1: Motivational Examples(a): Setting Tested Object Incorrectly**

Our approach's core innovations are twofold. First, **the explicit incorporation of project-aware knowledge provides the LLM with rich, contextualized information about the codebase**, enabling generation of more accurate and meaningful tests. Second, we introduce a **testing domain knowledge–guided separation of test design and test generation and multi-perspective test case design prompting strategy**. By decoupling the "what to test" from "how to test", we make the testing intent explicit and improve the semantic clarity and maintainability of the generated test code.

**Focal Method:**
```java
/** Puts a key-value mapping into Flat3Map<K, V>.
 * @param key   the key to add
 * @param value  the value to add
 * @return the value previously mapped to this key, null if none
 */
@Override
public V put(final K key, final V value) { ...}
```

**Test Code:**
```java
@Test
void testPut_NullKey_Size1() {
  flat3Map.put(null, "Value1");
  String oldValue = flat3Map.put(null, "Value2");
  assertNotNull(oldValue);        ❌ Inadequate Verification. The
                                     return value needs to be checked.
  // expected return is old value
  assertEquals("Value1", oldValue);
  // Verify size of map remains 1
  assertEquals(1, map.size());    ✅ Sufficient validation
                                     can detect errors.
}
```

**Figure 2: Motivational Examples(b): Insufficient Assert Statements**

We validate our approach on several real-world open-source projects, comparing it against state-of-the-art LLM-based baselines using both automated metrics (e.g., execution pass rate, line coverage) and a human study assessing readability, maintainability, and test intent clarity. Our findings show that (1) our method consistently outperforms baselines across six key metrics, improving execution pass rate by 5.69% and line coverage by 8.83% compared to the strongest baseline, while takes shorter time and generate less test cases; (2) the modular test case transformation component, which embodies the testing-principle-guided design, has the largest individual impact, with removal causing a drop of 24.08% in execution pass rate and 12.61% in line coverage; and (3) human evaluators rate our generated tests significantly higher in correctness, readability, and maintainability, confirming the practical advantages of our knowledge-driven framework.

Our main contributions are:

- We propose a novel test generation framework that integrates project-specific knowledge with testing domain knowledge. By leveraging comprehensive project knowledge and applying a separation of test design and test code generation guided by testing domain knowledge, our approach produces more accurate and maintainable test cases.
- We introduce a multi-perspective prompting strategy along with testing-domain-knowledge-guided separation of test case design and test method implementation, enabling the LLM to generate semantically rich, logically structured, and purposeful tests.
- We perform extensive evaluation, including automatic metrics and human studies, demonstrating that our method consistently outperforms existing LLM-based approaches in correctness, readability, maintainability, and overall quality.

All data/code used in this study is provided in the package [15].

## 2 Approach

We present KTester, a knowledge-aware unit test generation framework that leverages both project-specific knowledge and testing
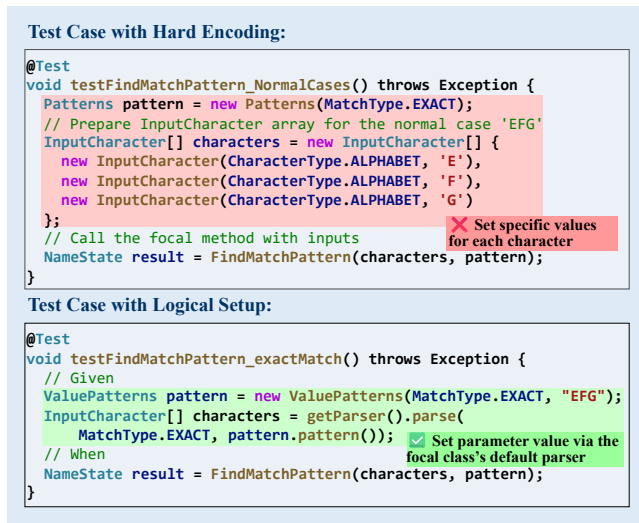
**Test Case with Hard Encoding:**

```java
@Test
void testFindMatchPattern_NormalCases() throws Exception {
  Patterns pattern = new Patterns(MatchType.EXACT);
  // Prepare InputCharacter array for the normal case 'EFG'
  InputCharacter[] characters = new InputCharacter[] {
    new InputCharacter(CharacterType.ALPHABET, 'E'),
    new InputCharacter(CharacterType.ALPHABET, 'F'),
    new InputCharacter(CharacterType.ALPHABET, 'G')
  };
  // Call the focal method with inputs          ✗ Set specific values
  NameState result = FindMatchPattern(characters, pattern);  for each character
}
```

**Test Case with Logical Setup:**

```java
@Test
void testFindMatchPattern_exactMatch() throws Exception {
  // Given
  ValuePatterns pattern = new ValuePatterns(MatchType.EXACT, "EFG");
  InputCharacter[] characters = getParser().parse(
    MatchType.EXACT, pattern.pattern());    ✓ Set parameter value via the
  // When                                     focal class's default parser
  NameState result = FindMatchPattern(characters, pattern);
}
```

**Figure 3: Motivational Examples(c): Hard Encoding Values**

domain knowledge to guide LLMs in generating high-quality, semantically meaningful, and maintainable unit tests.

Unlike prior LLM-based approaches that rely solely on the focal method as input and often generate brittle or unclear test cases, KTester decouples test design from test code generation, and injects knowledge at multiple stages to address key deficiencies of prior work (see Figure 1, 2 and 3). In particular, KTester is designed around the following core principles:

- Project Knowledge Awareness: We statically analyze the entire codebase to extract structured information about classes, methods, usage patterns, dependencies, and documentation. This knowledge is later used to resolve object instantiation, locate related APIs, and reduce hallucinations during code generation.
- Testing Knowledge Awareness: We explicitly guide the LLM to design test cases from multiple testing perspectives (e.g., control flow, boundary, exception handling), instead of treating the task as a black-box code-to-code translation.
- Modular Generation Pipeline: We separate test class framework construction, scenario-specific test design, and final test case synthesis into distinct stages, improving modularity, reusability, and interpretability of generated tests.

The two-stages process of KTester generating the test class is showed in Figure 4. First, during the offline knowledge extraction stage (Section 2.1), we perform static analysis of the target project to build a structured knowledge base capturing essential project-specific information such as class hierarchies, method signatures, field access patterns, call relationships, and document comments. This knowledge base provides a foundation for accurate, context-aware test generation. Second, in the online test generation stage (Section 2.2), given a focal method, KTester retrieves relevant project knowledge and testing heuristics to construct a knowledge-rich prompt for the LLM. The LLM then produces a structured test class framework along with diverse test cases, which are validated and

optionally repaired to ensure correctness and completeness. Figure 5 illustrates an example generated test class, including field declarations, setup and teardown methods, and multiple test methods.

## 2.1 Project Knowledge Extraction

To mitigate the limitations of LLMs in understanding project-specific semantics, KTester performs a comprehensive offline analysis to extract structured *project knowledge* from the codebase. This process is based on the assumption that unit tests are to be generated for methods within a project, and thus, knowledge beyond the focal method—such as class definitions, field usage, method invocations, and usage examples—is essential for meaningful test generation.

Specifically, KTester conducts two key analyses to construct a reusable project knowledge base: project structure knowledge mining, which captures the static architecture of the codebase—including class hierarchies, method signatures, and inter-method dependencies—and project usage knowledge mining, which extracts realistic invocation contexts that reveal how focal methods are constructed, initialized, and invoked in real code scenarios.

*2.1.1 Project Structure Knowledge Mining.* To enable knowledge-aware unit test generation, KTester statically analyzes the source code to mine structured *project structure knowledge*. This includes the architectural, semantic, and dependency-level information necessary for understanding the focal method's context and generating correct, maintainable test code. The analysis is based on abstract syntax tree (AST) parsing and code graph construction, allowing us to systematically extract and index key properties of classes and methods across the project.

Specifically, for each class and its methods or constructors, we extract the following information:

**Structural Metadata.** We collect the class name, package path, field declarations, constructor signatures, and method signatures (including parameter types and return types). This information enables the model to understand how to instantiate the target class or construct inputs for the method under test. For example, if a method requires a parameter of type UserConfig, this metadata reveals how UserConfig can be constructed via its fields and constructors.

**Document Comments.** We extract Javadoc-style comments associated with classes, constructors, and methods, as they often contain high-level semantic descriptions of behavior and usage constraints. These comments help LLMs infer testing intent and expected behavior. For instance, a comment like /* Returns null if no user is found */ suggests that null should be included as a valid return value in test oracles.

**Dependency Relations.** For each method or constructor, we extract the list of invoked methods and accessed fields. These relationships are fundamental for understanding how the method interacts with its internal state or other classes. For example, if method updateBalance() internally calls validateAccount() and accesses this.balance, such dependencies indicate that tests should ensure the account is in a valid state and that balance updates are properly asserted. These dependency relations also support the identification of functionally related methods (used later in Usage
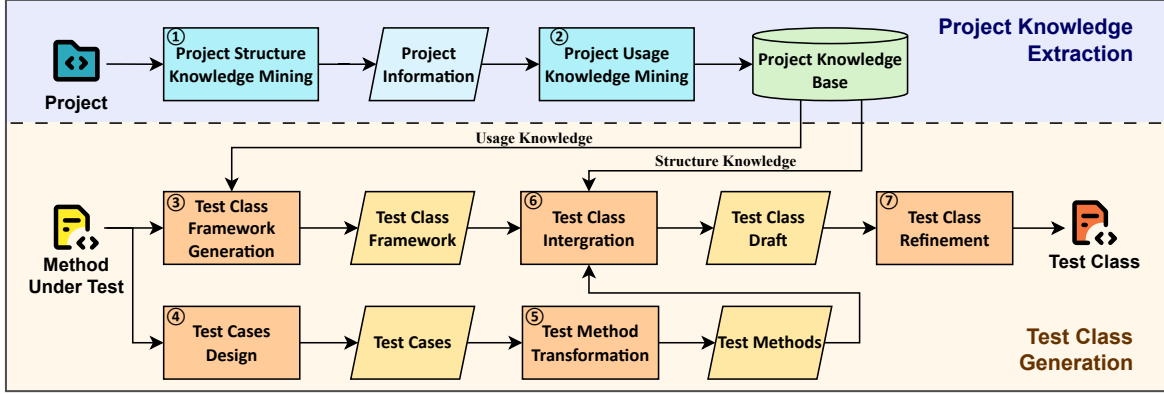
**Figure 4: The Framework of KTester**

Context Construction) and enable accurate call path tracing for generating invocation examples (Section 2.1.2).



**Figure 5: Test Class Example**

*2.1.2 Project Usage Knowledge Mining.* To mitigate the limitations of LLMs in understanding project-specific usage patterns, we extract *Project Usage Knowledge* from the source code. This knowledge captures realistic invocation scenarios of the focal method, including (1) how its input parameters are constructed and (2) how the target object is initialized—both of which are essential for generating executable and semantically meaningful unit tests.

Such usage knowledge provides concrete, in-project examples that reflect actual calling contexts, enabling the LLM to synthesize inputs that align with real usage scenarios. For example, it reveals how the focal class is typically constructed in the project, which may involve complex initialization logic or dependency injection. It also uncovers how input parameters are derived—such as being composed from helper methods, shared resources, or intermediate computations—rather than being hardcoded literals. By grounding the generation process in real-world usage, this approach improves both the realism and maintainability of the generated test cases. As



**Figure 6: Illustration of extracting invocation examples from a caller method via control flow pruning.**

the code snippet illustrated in Figure 6, the focal method *findMatchPattern(InputCharactor[], Patterns)* is invoked inside *findMatchPattern(ValuePatterns)* after checking whether the variable *pattern* is an instance of *ValuePatterns* class, and called a Parser to parse this pattern into an *InputCharactor* array. Such code would be extracted as representative usage knowledge.

To mine project usage knowledge, we begin by identifying the relevant caller context and then extract semantically meaningful execution traces that illustrate how the focal method is invoked in practice.

**Caller Method Discovery.** For each focal method, we identify its *caller methods* using the static call graph built during the offline analysis. When the focal method is private or not directly reachable from public APIs, we trace the shortest call chain from an externally accessible method to the focal method. This ensures that the usage knowledge reflects a realistic and complete setup process, capturing how the focal method is constructed and triggered in actual usage.

**Path-sensitive Usage Trace Extraction.** For each discovered caller method, we construct a Control Flow Graph (CFG) [16] and perform intra-procedural slicing to extract execution paths from the method's entry point to the focal method invocation. By pruning unrelated branches and removing irrelevant operations, we isolate minimal yet semantically rich usage traces. These traces, which

demonstrate how the focal class is instantiated and how arguments are prepared, are incorporated into the generation prompt as part of the Project Usage Knowledge. This provides the LLM with grounded, context-aware examples that help it synthesize test inputs aligned with the project's actual usage conventions.

*2.1.3 Function-Level Index Construction.* To enable scalable and consistent test generation, each method or constructor is transformed into a *knowledge unit* containing its signature, documentation, and dependencies (invoked methods and accessed fields). These units are stored in an indexable format, forming a function-level knowledge base that supports efficient retrieval of related methods, usage contexts, and class-level information during prompt construction.

Unlike ad hoc, per-method extraction, our analysis is conducted *once per project* and reused across focal methods. The knowledge base is also *incrementally updatable*—only changed entities need to be reprocessed—making it practical for real-world, evolving codebases. While our method is language-agnostic, we focus on Java due to its popularity and mature analysis tooling[17, 18]. We leverage Spoon [18] to extract ASTs, control flow, and dependency graphs efficiently.

## 2.2 Test Class Generation

As illustrated in Figure 4, KTester generates a complete and executable test class for a given focal method by leveraging offline-extracted project knowledge and domain-specific testing expertise. It constructs a rich generation context that integrates structural program information, usage patterns, and unit testing heuristics, guiding the LLM to produce high-quality, realistic unit tests. The generation pipeline consists of five steps: test class framework generation, test case design, test method transformation, test class integration, and test class refinement.

Inspired by how developers and testers write unit tests in practice, KTester mimics the typical testing workflow—first planning test cases based on expected behaviors and coverage goals, then gradually transforming them into executable test code—rather than relying on a single-step code generation process.

*2.2.1 Test Class Framework Generation.* The first step prompts the LLM to generate a test class framework that establishes the necessary environment for unit testing. Unlike prior methods [19, 20] that directly prompt the LLM with a focal method and expect complete test cases in one shot, KTester decouples the generation process by first constructing a reusable, project-aware test class framework. This decision is grounded in a key insight: test generation is not merely a method-level translation task, but a context-sensitive activity requiring awareness of surrounding project structure, usage conventions, and domain-specific testing patterns. As illustrated in Figure 7, we design a prompt template grounded in project knowledge and test domain knowledge to guide this process. The prompt is composed of several key components, including core task and instructions, focal method and class, relevant project knowledge, test class template, and output specification.

The **core task and instructions** section provides high-level guidance to the LLM on setting up proper initialization and cleanup routines, enforcing access restrictions to target class internals, and

ensuring syntactic correctness, thereby enabling a robust and maintainable test framework. The **focal method and class** section specifies the method under test and its containing class, providing the primary scope for framework generation. The **output specification** section defines the expected format and content of the generated framework, ensuring consistency and completeness.

To enrich the context, the **relevant project knowledge** section of the prompt draws from an offline-constructed knowledge base that captures both semantic and structural details of the codebase. This includes document comments for the focal class and method, which describe functionality, preconditions, and side effects to inform meaningful setup and teardown procedures. It also comprises constructors and parameter details, assisting the LLM in object instantiation and input preparation, particularly in complex scenarios. Additionally, usage knowledge mined from the codebase provides real-world usage patterns of the focal method and related classes, grounding the test framework in practical contexts. Finally, existing test classes associated with the focal class are incorporated as references for mocking strategies, resource management, and testing conventions, with only setup logic and class-level configurations extracted to avoid redundancy.

The **test class template** section outlines the structural skeleton, including annotations, lifecycle methods (e.g., setup and teardown), and placeholders for field declarations. This template, informed by practical experience in writing test code, serves as a starting point for the LLM to fill in with project-specific logic and configurations.

By integrating these components, KTester creates context-aware prompts that enable the LLM to generate reusable, well-structured test class frameworks aligned with project conventions. This modular approach enhances consistency and reduces errors common in one-shot generation. Separating framework from test logic also enables reuse across methods, minimizing redundant effort and reflecting practical testing workflows.

*2.2.2 Test Cases Design.* After generating the test class framework, the next step is to design high-quality test cases that cover different functional and non-functional behaviors of the focal method. In this phase, KTester focuses on the planning of test cases, rather than directly generating test code. Unlike the previous step, test case design solely leverages domain-specific testing knowledge, without relying on project-specific context. The only input is the focal method itself.

KTester adopts a multi-view guidance strategy grounded in established software testing principles. The LLM is explicitly instructed to design test cases from multiple perspectives—such as functional behavior, boundary conditions, and exception handling—and to organize them into logical groups based on shared testing intent. Each test case is represented internally in a machine-readable intermediate format that records structured natural-language descriptions of the scenario. In our implementation, this intermediate representation is realized using a lightweight JSON structure, and an illustrative example is provided in our replication package [15]. This design intentionally decouples test intent from concrete executable code, improving clarity and maintainability. The intermediate representations are subsequently translated into framework-specific test code.

---

**Prompt For Test Class Framework Generation**

You are tasked with **creating a framework for a test class in Java**. Now analyze the following information, and **complete the test class template** in the following requirements:

1. Provide the code wrapped in a Markdown code block "```java" at the beginning and "```" at the end.
2. Initialize necessary fields and dependencies within *setupBeforeAll()* and *setupBeforeEach()*.
3. Handle any required cleanup in *teardownAfterEach()* and *teardownAfterAll()*.
4. ...

<div align="right"><b>Task Description</b></div>

**@input{focal method and target class}**
Here is the code of focal method, *<method sigature>*, with the focal class *<class name>*:

```
Class <class_name> {
    dependent field declaration;
    dependent method signature;
    <method signature> {<code>} }
```

<div align="right"><b>Code Under Test</b></div>

**@input{related context}**
Document comment of focal class & focal method: *<javadoc>*
Constructors of focal class: *<constructor code>*
Constructors of param *<param class>*: *<constructor code>*
Usage examples of focal method: *<usage example code>*
Existing Test class: *<example test code>*

<div align="right"><b>Relevant Project Knowledge</b></div>

**@input{test class template}**

```
package <package name>;
import ...;
// Add other imports if necessary
Class <class_name> {
    @BeforeAll
    static void setupBeforeAll() {}
    @BeforeEach
    void setupBeforeEach() {}
    @AfterEach
    void teardownAfterEach() {}
    @AfterAll
    static void teardownAfterAll() {} }
```

<div align="right"><b>Test Class Template</b></div>

**Figure 7: Prompt for test class framework generation.**

To facilitate structured generation, the LLM is required to output grouped test cases, where each group targets a distinct aspect of the focal method. For example, one group may cover valid inputs for core functionality, another may focus on boundary values such as empty lists or null, and a third may address failure or exception-triggering scenarios. This grouped organization improves modularity and controllability during subsequent transformation steps and mirrors human testing practices, where related scenarios are often clustered for clarity and reuse.

To guide the LLM, we design three types of prompts, each corresponding to a classical testing perspective, as shown in Figure 8. These prompts share a common structure, with variations only in the design guidance section.

- **Condition branch-driven prompt**: This prompt encourages the LLM to explore different control-flow paths, such as conditional statements and loops, improving branch and path coverage.
- **Functionality-driven prompt**: This prompt focuses on the intended semantics of the method, helping the LLM generate test inputs that exercise typical behaviors and edge cases derived from parameter types and return values.
- **Exception-oriented prompt**: This prompt drive the generation of test cases for robustness and failure handling,

**Prompt For Test Cases Design**

As a professional Java software testing expert, your task is to **generate formatted test cases for a focal method**.
Following the guidance: *<conditional branch> or <functional behavoir> or <exception behavior>*

| Conditional Branch Guidance | Functional Behavior Guidance | Exception-oriented Guidance |
|---|---|---|
| Create test cases to cover all **conditional branches** in the target method, including *if, else if, else, switch* and *case* statements. | Analyze the **method's functionality** and identify possible input combinations including **normal cases** and **boundary conditions**. | Analyze method behavior and identify ALL **possible exception scenarios**, then design test cases that trigger each identified exception. |

**@input{focal method and target class}**
Here is the code of focal method, *<method sigature>*, with the focal class *<class name>*: *<code under test>*

**@output{test cases}**
Follow the JSON format provided here for the output:

```
[{
    "group": "test name",
    "cases": [{
        "input": [
            {"parameter": "param name", "value": "param value"},
            {...}],
        "expected": "expected exception or behavior",
        "description": "test scenario description"
    }, {...}]
},{...}]
```

**Figure 8: Prompts for multi-view test case design.**

including invalid or unexpected inputs that should trigger exceptions or error states.

This multi-view prompt design is modular and extensible—new testing views (e.g., performance, concurrency) can be incorporated by adding corresponding prompt templates. By separating test planning from execution and grounding the design in testing theory, this step enhances the quality, coverage, and interpretability of generated tests.

*2.2.3 Test Method Transformation.* After designing structured test cases, KTester transforms each test case group into an executable test method, aligning with how developers typically write tests—first decide what to test, then implement the logic.

**Prompt Construction.** We design a usage-aware prompt with several semantically distinct sections to guide the LLM in generating executable test methods, as shown in Figure 9. First, test case groups contain structured cases sharing a common testing intent (e.g., valid inputs, edge cases, exceptions), each with a scenario, inputs, and expected output. The LLM generates one test method per group to ensure modularity and clarity. Second, the usage context, retrieved from the project knowledge base (Section 2.1.2), provides focal method dependencies and related methods to guide realistic invocation and assertion generation. Lastly, the test class framework offers reusable setup/teardown code and field declarations, promoting consistency and eliminating redundancy.

**Usage Context Retrieval.** To enhance generation quality, we retrieve methods that are functionally related to the focal method and its dependencies, offering a broader context of real-world usage patterns. Based on our observations, related methods are identified based on shared usage of functions and fields. We define the similarity between two functions $a$ and $b$ using a Jaccard-based

---

**Prompt For Test Method Transformation**

You are a professional Java software testing expert. Your task is to **generate unit test method** for a focal method, based on existing test class framework and formatted test cases. Requirements are as follows:
1. Follow the given test class framework; do not modify *@Before/@After* methods.
2. Create **one test function per test case group**, ensuring all cases are included. Use parameterized tests where appropriate.
3. ...

**Task Description**

**@input{focal method and target class}**
Here is the code of focal method, *<method sigature>*, with the focal class *<class name>*: *<code under test>*

**Code Under Test**

**@input{initial test class framework}**
                    *<test class framework code>*
**@input{existing test case}**
Here's the existing test cases: *<json test cases>*

**Test Class Framework & Formatted Test Cases**

**@input{related context}**
Document comment of focal class & focal method: *<javadoc>*
Dependent classes of focal method:

```
Class <class_name>:
    relationship: param/return type /related class
    accessed field: <field declaration>
    invoked method: <javadoc> + <method signature>
    related function: <function signature>
                related with <method signature>
```

**Project Knowledge Context**

**Figure 9: Prompt for test method transformation.**

metric [21] by combining with method usage similarity $Sim_m(a, b)$ and field usage similarity $Sim_f(a, b)$:

$$Sim(a, b) = Sim_m(a, b) + Sim_f(a, b) \quad (1)$$

$$Sim_m(a, b) = \frac{|M(a) \cap M(b)|}{|M(a) \cup M(b)|} \quad (2)$$

$$Sim_f(a, b) = \frac{|F(a) \cap F(b)|}{|F(a) \cup F(b)|} \quad (3)$$

Here, $M(a)$ and $F(a)$ denote the sets of methods and fields used by function $a$, respectively. We select the top-N most similar functions to construct the context.

For each related class, we include: (1) the class declaration, optionally with Javadoc comments; (2) relevant field declarations accessed by dependent or related methods; (3) signatures of dependent methods with comments; and (4) signatures of related methods, annotated with markers indicating shared usage. This context provides the LLM with concise, relevant information to support realistic test generation.

*2.2.4  Test Class Integration.* After transforming individual test cases into test methods, the next step is to integrate them into a coherent and executable test class. This step ensures that all test methods, including newly generated ones and any accompanying setup or teardown logic, are correctly incorporated into the test class framework without introducing redundancy or conflicts.

To achieve this, KTester performs static analysis on the generated methods to identify their roles and resolve duplicates. The integration process proceeds as follows:

- **Setup/Teardown Merging.** If a newly generated method is identified (via static analysis) as a setup or teardown method

(e.g., annotated with @Before, @BeforeEach, etc.), and an existing method of the same type already exists, we merge their bodies. Specifically, any code in the new method that is not present in the original is appended, ensuring initialization routines are preserved and extended without duplication.
- **Test Method Deduplication.** For test methods and general helper methods, we identify duplicates by comparing their method names. When two methods share the same name, we further compare their bodies and retain the version with the longer implementation, using length as a proxy for logical richness and completeness. This deduplication addresses a common pattern where the LLM initially produces simple placeholder tests and later generates more complete versions with the same or similar names, ensuring that the final output preserves the more comprehensive implementation.
- **Non-conflicting Insertion.** For all other non-conflicting methods, we append them directly to the test class. This modular addition supports test diversity and extensibility without disrupting the original structure.

The integration process leverages standard AST-based static analysis tools to parse and manipulate the Java code structure safely and consistently. This automation ensures that the final test class is executable, maintainable, and free of duplicate or conflicting code.

By decoupling method generation from class integration, our approach maintains modularity while aligning with practical software engineering workflows, where tests are often extended iteratively and merged across sessions or contributors.

*2.2.5  Test Class Refinement.* In this stage, KTester validates the generated test class by checking for compilation and execution errors. Despite the rich contextual information used during generation, LLM outputs may still contain issues such as hallucinated import statements, improper mocking, incorrect usage of private methods, or runtime errors during test execution. To address these problems, we adopt an iterative refinement process consisting of both static validation and dynamic repair.

**Rule-based Repair.** Rule-based repair focuses on correcting straightforward and common errors detected by static analysis, such as hallucinated or missing import statements and unresolved symbols. Leveraging the previously constructed project-level test knowledge base, KTester can automatically resolve dependencies by appending missing imports or correcting incorrect ones based on simple class names. The comprehensive indexing of classes, methods, and fields from the static analysis phase facilitates accurate dependency resolution without human intervention.

**LLM-based Repair.** For errors arising from semantic misunderstandings or subtle logic issues that cannot be fixed by static rules alone, KTester uses LLM-based repair. This involves extracting detailed compile-time error messages, runtime error logs, and failing test feedback. This information, combined with the focal method context and relevant function signatures, is incorporated into a repair prompt that is fed back to the LLM. The model then performs targeted correction to resolve issues such as incorrect assertions, missing exception handling, or logic bugs causing runtime failures.

Importantly, this refinement process extends beyond compilation to include execution feedback. If runtime errors occur or tests fail

due to assertion mismatches, these failure details are fed into the repair pipeline to enable iterative correction and improvement of the generated test code. This dynamic feedback loop helps ensure that the final test class is not only syntactically correct but also functionally robust and reliable.

## 3 Evaluation

In this section, we evaluate the effectiveness of KTester by answering the following research questions (RQs):

**RQ1 (Effectiveness Comparison)**: Does KTester outperform state-of-the-art baselines in terms of coverage scores and execution pass rate when testing complex methods?

**RQ2 (Ablation Study)**: How do the key components of KTester contribute to its overall performance?

**RQ3 (User Study)**: Are the tests generated by KTester more readable and maintainable than those produced by baseline methods?

### 3.1 Experimental Setup

*3.1.1 Dataset.* To evaluate the effectiveness of KTester, we adopt the HITS dataset [20], which is constructed from 10 popular open-source Java projects spanning various domains (e.g., microservices, command-line tools, event engines). This dataset specifically targets complex methods—defined as those with cyclomatic complexity greater than 10—making it well-suited for assessing how effectively LLMs handle methods with intricate control flow and dependencies. The dataset comprises 110 tasks, each consisting of a focal method and its containing class.

**Table 1: Detailed Information of HITS Dataset**

| Project | Domain | Version | #MUT |
|---|---|---|---|
| Commons-CLI | Cmd-line Interface | 1.7.0-SNAPSHOT | 2 |
| Commons-CSV | Data processing | 1.10.0 | 6 |
| Gson | Serilization | 2.10.1 | 20 |
| Commons-codec | Encoding | 3a6873e | 18 |
| Commons-collections4 | Utility | 4.5.0-M1 | 14 |
| JDom2 | Text Processing(XML) | 2.0.6 | 21 |
| Datafaker | Data Generation | 1.9.0 | 6 |
| Event-ruler | Event Engine | 1.4.0 | 15 |
| windward | Micoservices | 1.5.1-SNAPSHOT | 2 |
| batch-processing-gateway | Cloud Computing | 1.1 | 6 |

*3.1.2 Baselines.* We compare KTester with 4 state-of-the-art test generation methods: three are purely LLM-based, and one integrates search-based software testing (SBST) with LLMs. The LLM-based baslines adopt distinct strategies for context construction and test generation:

**ChatUnitTest** [19] provides the focal method and handcrafted context to the LLM, and applies iterative repair using error feedback from failed executions.

**ChatTester** [22] is similar to ChatUnitTest but differs in how it constructs the focal method's context, building it incrementally as needed.

**HITS** [20] decomposes the focal method into slices and prompts the LLM to generate test cases for each slice, which are then merged into a complete test suite. This fine-grained approach helps mitigate the difficulty LLMs face when reasoning about complex logic holistically.

For fair comparison, we adapt official or open-source implementations of these methods [23]. ChatUnitTest and ChatTester generate one test class per focal method, while HITS's slice-level tests are merged into a single class to ensure consistent execution and setup. All baselines and KTester use the same LLM backend—gpt-4o-mini [24]—with default API settings with temperature=1.0, ensuring consistency in generation quality, efficiency, and cost.

The SBST–LLM hybrid method, **UTGen** [25], employs a traditional SBST tool to generate initial test cases, then leverages an LLM to rename functions and variables in order to improve test readability. Since UTGen rely on an LLM's ability on code understanding instead of enhancing code coverage, we follow the initial configuration provided in its replication package [26] (EvoSuite [27] as the generation tool and CodeLlama-7b [28] used for refinement) when conducting our experiments to evaluate its performance.

### 3.2 RQ1: Effectiveness Comparison

*3.2.1 Design.* We use the HITS dataset (Section 3.1.1), which contains 110 test generation tasks from 10 open-source Java projects. Our method builds the project knowledge base by analyzing the corresponding project source code versions in the dataset. For these tasks, both KTester and baseline methods generate unit tests using GPT-4o-mini as the backbone model. For fair comparison, KTester and all baselines generate exactly one test class per focal method, with no restriction on the number of test cases or generation time within each class. This allows complex methods to achieve high coverage through diverse test cases. We limited automatic repair iterations to 5 in KTester. To mitigate the effect of randomness introduced by LLMs, we repeated each experiment three times and report the average results.

To systematically evaluate the effectiveness of generated unit tests, we adopt eight widely used metrics for unit test quality, capturing correctness, sufficiency and efficiency [19, 20, 22, 29].

- **Compile Pass Rate (CPR)**: percentage of test classes that compile successfully.
- **Execution Pass Rate (EPR)**: percentage that run without runtime errors or assertion failures.
- **Line Coverage (LC)**: ratio of executable lines in the focal method covered by all generated tests.
- **Branch Coverage (BC)**: percentage of conditional branches covered.
- **Line Coverage of Passed Tests (LCP)**: line coverage computed only on tests that compile and execute successfully.
- **Branch Coverage of Passed Tests (BCP)**: branch coverage computed similarly on passing tests.
- **Average Time pre Task (AvT)**: average time required by a method to generate test class for a single task.
- **Average Test Cases Number (AvTC)**: average number of test cases contained in each generated test class.

Coverage is measured using Jacoco [30], an open-source tool for measuring code coverage in Java projects, ensuring objective and

**Table 2: Comparison of Effectiveness across Baselines**

|        | UTGen  | ChatTester | ChatUniTest | HITS   | KTester |
|--------|--------|------------|-------------|--------|---------|
| CPR    | **100** | 55.18     | 98.97       | 97.82  | **100** |
| EPR    | **90.05** | 50.17   | 65.26       | 71.38  | 77.07   |
| LC     | 31.17  | 28.51      | 44.64       | 52.27  | **61.10** |
| BC     | 28.72  | 24.94      | 38.06       | 45.93  | **52.59** |
| LCP    | 30.69  | 22.26      | 33.47       | 43.85  | **54.49** |
| BCP    | 28.11  | 19.63      | 28.51       | 38.81  | **46.23** |
| AvT (s)| 1068   | 354.83     | 200.96      | 625.69 | **152.68** |
| AvTC   | 9.23   | 3.32       | 4.77        | 15.78  | 7.33    |

consistent evaluation. By combining correctness (CPR, EPR), coverage metrics (LC, BC, LCP, BCP) and efficiency metrics (AvT, AvTC), we provide a balanced and robust assessment of test quality, facilitating fair comparison between KTester and baseline approaches.

*3.2.2 Results.* Table 2 presents the effectiveness comparison across four methods. Overall, KTester achieves the best performance on all six evaluation metrics, demonstrating its ability to generate high-quality unit tests both in terms of correctness and coverage. **Correctness.** KTester reaches a perfect Compile Pass Rate (CPR) of **100%**, indicating that all generated test classes are syntactically valid and type-correct. While UTGen also achieves 100% CPR, the other LLM-based baselines fall short, suggesting potential issues in code generation or dependency handling. Regarding the Execution Pass Rate (EPR), which measures runtime correctness, KTester attains 90.05%, slightly lower than UTGen. This is expected, as UTGen builds upon EvoSuite-generated tests, which provide a strong foundation for runtime correctness. Nevertheless, KTester still surpasses all other LLM-based methods in EPR, demonstrating its effectiveness in generating executable and correct test cases. This result indicates that incorporating project knowledge yields substantial benefits, as it more effectively guides LLMs toward generating executable test code.
**Sufficiency.** In terms of code coverage, KTester achieves the highest Line Coverage (LC) of **62.78%** and Branch Coverage (BC) of **54.71%**, significantly outperforming HITS (49.74% LC, 43.74% BC) while generating less test case (7.33 vs. 15.78). ChatTester yields the lowest coverage, highlighting its limited ability to explore program logic. UTGen is constrained by EvoSuite's inability to operate on certain project-specific packages (e.g., *com.apple.spark* in the batch-processing-gateway project) and by the path-explosion problem inherent to SBST techniques. As a result, its coverage exceeds only that of ChatTester. When restricting evaluation to only the tests that both compile and execute successfully, KTester maintains the lead with **54.49%** Line Coverage of Passed Tests (LCP) and **46.23%** Branch Coverage of Passed Tests (BCP), confirming the practical adequacy of its outputs.

These results confirm that integrating project-specific and test-domain knowledge into LLM-based generation pipelines can significantly enhance the validity and adequacy of produced test cases. In particular, KTester demonstrates not only higher syntactic and semantic correctness, but also stronger capability in exercising program logic and control flow, making it a more effective and practical solution for automated unit test generation.

**Test Class (generated by HITS)**

```java
public class Flat3Map_2_0_Test {
  @Test
  void testGet_MatchingKey_ReturnsValue() {
    Flat3Map<String, Integer> flat3Map=new Flat3Map<>();
    Field sizeField=Flat3Map.class.getDeclaredField("size");
    sizeField.setAccessible(true);
    sizeField.set(flat3Map, 1);    Using reflection to access
    ...                            private fields
    Integer expectedValue = 100;
    Integer actualValue = flat3Map.get("key1");
    assertEquals(expectedValue, actualValue);
}}
```

**Test Class (generated by KTester)**

```java
class Flat3Map_get_Test {
    private Flat3Map<Integer, String> flat3Map;
    @BeforeEach
    void setupBeforeEach() {
        flat3Map = new Flat3Map<>();
        flat3Map.put(1, "One");
        flat3Map.put(2, "Two");     Using related method to
        flat3Map.put(3, "Three");   set up
  } @Test
    void testGet_ExistingKey() {
        assertEquals("One", flat3Map.get(1), "...");
        assertEquals("Two", flat3Map.get(2), "...");
        assertEquals("Three", flat3Map.get(3), "...");
}}
```

**Figure 10: Test Classes generated by HITS and KTester.**

**Table 3: Comparison of Effectiveness across KTester implemented with different models**

|       | KTester-gpt | KTester-claude | KTester-deepseek |
|-------|-------------|----------------|------------------|
| CPR   | 100         | 100            | 93.92            |
| EPR   | 77.07       | 81.75          | 82.41            |
| LC    | 61.10       | 66.46          | 71.22            |
| BC    | 52.59       | 60.66          | 65.75            |
| LCP   | 54.49       | 56.30          | 65.19            |
| BCP   | 46.23       | 50.76          | 59.33            |

Figure 10 shows the test class generated for the *get()* method with HITS and KTester. HITS used reflection to access and set private values, which increases the difficulty of understanding. In fact, since Flat3Map is a Map data structure, we can better modify private fields through the relevant methods *put ()* in this class, thus covering every conditional branch during testing.
**Generalizability.** To assess whether KTester maintains its performance when implemented with alternative LLMs, we implemented KTester using claude-3.5-haiku-20241022 [31] and deepseek-v3.1 [32], and the experimental results are reported in Table 3, which achieved even better results than gpt4o-mini, confirming its generalizability.

**Finding 1:** KTester outperforms all baselines across eight metrics. Compared to the strongest baseline (HITS), and while generating on average 8.45 fewer test cases, it improves execution pass rate by 5.69% and line coverage by 8.83%, demonstrating clear gains in both correctness and sufficiency.

**Table 4: Comparison of Effectiveness across Variants and original KTester.**

|  | CPR | EPR | LC | BC | LCP | BCP |
|---|---|---|---|---|---|---|
| KTester-UTE | 100 | 58.20 | 56.84 | 50.48 | 44.54 | 38.90 |
| KTester-FMR | 98.15 | 53.91 | 58.19 | 50.95 | 45.40 | 39.41 |
| KTester-MVG | 100 | 64.57 | 54.99 | 48.59 | 43.80 | 38.89 |
| KTester-DGT | 97.76 | 52.99 | 48.49 | 44.05 | 35.79 | 31.99 |
| KTester | **100** | **77.07** | **61.10** | **52.49** | **54.49** | **46.23** |

## 3.3 RQ2: Ablation Study

*3.3.1 Design.* To assess the individual impact of key components in KTester on unit test generation, we conducted an ablation study. Four variants were created, each removing or modifying one component while keeping the rest of the pipeline intact. This design isolates each component's effect on correctness and coverage. The variants include:

- **KTester-UTE**: without the procedure of **u**sage **t**race **e**xtraction (section 2.1).
- **KTester-FMR**: without the step of **f**unctionally related **m**ethod **r**etrieval (section 2.1.2).
- **KTester-MVG**: using only conditional branch guidance instead of the whole **m**ulti-**v**iew **g**uidance strategy (section 2.2.2).
- **KTester-DGT**: replace the processes of test cases design and test method transformation (section 2.2.3) with a single process that **d**irectly **g**enerates **t**est methods using a multi-view guidance strategy.

Using the same experimental setup and metrics as in RQ1, we evaluated all variants and compared their results to the original KTester to measure each component's contribution to test quality.

*3.3.2 Results.* Table 4 shows the performance of each KTester variant compared to the original method. All variants experience a drop across correctness and coverage metrics, confirming that each key component positively contributes to overall test quality. Among them, KTester-DGT, which removes the intermediate structured test case generation step, causes the most significant degradation. For example, Execution Pass Rate (EPR) drops by about 24.08% (from 77.07% to 52.99%) and Line Coverage (LC) decreases by approximately 12.61% (from 61.10% to 48.49%). This indicates that separating test case design from test method generation is crucial for effectiveness. The other variants—KTester-UTE, KTester-FMR, and KTester-MVG—also reduce performance, but their impacts are less distinct and harder to differentiate clearly. This suggests that while invocation pattern extraction, similar function retrieval, and multi-view prompting all improve results, the modular test case transformation step plays the most critical role.

**Finding 2:** The modular test case transformation (KTester-TCG) has the largest impact on performance, with EPR and LC dropping 22.15% and 14.29%, respectively. Other components also contribute positively but with less distinguishable effects.

## 3.4 RQ3: User Study

*3.4.1 Design.* To investigate whether the tests generated by our approach exhibit superior readability and maintainability, we conducted a user study targeting professional developers. We first selected 10 tasks from the dataset where both KTester and the baseline methods successfully produced compilable test cases. Ensuring compilability allows participants to concentrate on qualitative properties of the test code—such as clarity, structure, and maintainability—without being influenced by technical correctness issues.

We recruited participants through a public invitation distributed across the computer science departments of six universities. From the volunteers, we selected 15 participants (5 Ph.D. students and 10 Master's students), all with 2–5 years of Java development experience. Participants received compensation to encourage careful and unbiased evaluation.
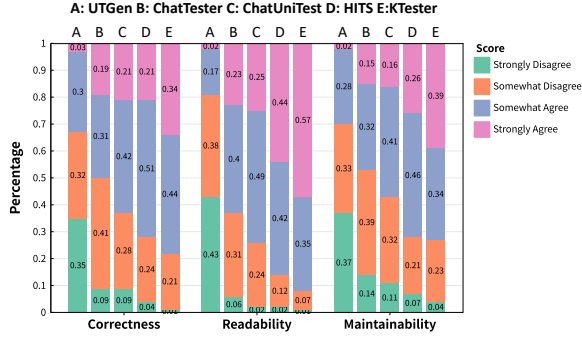
For each task, participants were presented with the focal method under test, its API documentation and relevant class context, then the related test classes generated by different methods. To avoid bias, the identities of the methods were fully anonymised, and the presentation order was randomised. Participants were instructed to respond to the following three questions, each designed to target a key quality dimension of the generated test code:

(1) To what extent do you agree that the unit test class is functionally correct? Specifically, does it adequately cover relevant input combinations, contain at least one meaningful assertion per test method, and reliably detect functional faults through deterministic and repeatable execution?
(2) To what extent do you agree that the unit test class is highly readable? For instance, are the method and variable names clear and descriptive? Does each test follow the Arrange–Act–Assert structure to clearly communicate intent? Are assertion failure messages informative and helpful for debugging?
(3) To what extent do you agree that the unit test class is maintainable? For example, does it minimise redundancy or duplicated code? Is the test class support the easy addition of new test cases? Can developers easily update affected tests when the source code changes?

Each question is rated on a four-point Likert scale—strongly agree, somewhat agree, somewhat disagree, and strongly disagree. This setup enables quantitative comparison of the subjective quality of tests generated by different approaches and allows us to assess whether our method consistently produces more developer-friendly outputs. Following prior work [22, 33], we adopt an even-numbered scale to avoid ambiguous "neutral" responses; specifically, the four-point design clearly separates two positive options from two negative ones.

*3.4.2 Results.* Figure 11 presents the distribution of developer ratings across the three evaluation dimensions. Across all three dimensions, our method (KTester, labeled as E) consistently outperforms the baselines in terms of perceived quality.

In the dimension of Correctness, KTester achieves the highest proportion of "Strong Agree" ratings (0.34), while keeping "strongly disagree" scores low (0.01). This trend indicates a strong alignment

**Figure 11: Score Distribution of Correctness, Readability and Maintainability**

with expected test behavior as perceived by developers. In terms of Readability, KTester demonstrates the most favorable distribution, with over 90% of responses rated as "strongly Agree" (0.57) or "Somewhat Agree" (0.35), suggesting that participants found its tests clearer and more accessible compared to those from other methods. For Maintainability, KTester has the highest proportion of "Strong Agree" responses (0.39) and the lowest proportion of "strongly disagree" scores (0.04). This suggests that developers found the structure and logic of the tests generated by our approach easier to understand and adapt.

Interestingly, although UTGen leverages a LLM to improve the readability of test code, it still received more than 67% negative responses across all three dimensions. We attribute this outcome to two main factors: (1) In many cases, some function signatures and variable names are not modified by the LLM, leaving the test code difficult for participants to understand. (2) EvoSuite invokes private methods through public caller methods, whereas LLM-based methods rely on reflection. When the source code changes, participants found it difficult to determine the scope of impact for the test cases produced by UTGen, which negatively affected their assessments of correctness and maintainability.

> **Finding 3:** KTester consistently outperforms all baselines in human evaluation across correctness, readability, and maintainability.

## 4 Related Work

In this section, we mainly introduce related work on LLM-based unit test generation.

### 4.1 LLM-based Unit Test Generation

Large language models (LLMs) have become a powerful tool for automated unit test generation, offering human-like test code that overcomes the readability and maintainability limitations of traditional methods [34–38]. Existing LLM-based approaches can be broadly categorized into two paradigms: fine-tuning and prompt-based methods.

**Fine-tuning** approaches treat test generation as a supervised learning task, typically formulated as a sequence-to-sequence translation problem. A representative example is AthenaTest [39], which

fine-tunes BART models on curated datasets where focal methods are inputs and complete test cases are outputs. These models learn explicit mappings between source code and corresponding tests using domain-specific corpora. Variants of this paradigm [40–42], including TeCo [43], explore different encoder-decoder architectures to better capture the semantic relationships between code and the associated test logic. A recent example is EXLONG [44], fine-tuned from CodeLlama to generate exceptional behavior tests. Unlike such fine-tuning approaches, our method injects project and testing knowledge via prompt design and generation workflow, making it more flexible and adaptable to stronger future LLMs.

**Prompt-based** methods, in contrast, rely on the zero-shot reasoning abilities of instruction-tuned models without additional training. Systems such as ChatTester [22], ChatUniTest [19] and ASTER [33] use carefully designed prompts to guide models like GPT in generating test suites. These methods often adopt a multi-step reasoning process, decomposing tasks into subcomponents such as test identification, input generation, and assertion formulation. TestPilot [45] enhances this strategy by incorporating documentation and usage examples, while TestSpark [46] combines prompt-based generation with search-based techniques in IDEs.

Recent work explores hybrid approaches that integrate multiple strategies. CODAMOSA [47] combines LLM-generated test seeds with evolutionary algorithms to improve coverage-driven test generation. CoverUp [48] employs iterative dialogues to refine outputs based on coverage feedback. SymPrompt [49] leverages symbolic execution to guide prompting.

Despite these advancements, major challenges remain. Most methods lack access to project-specific context, such as class instantiation, module dependencies, and API usage patterns—leading to syntactically correct but semantically invalid tests. Moreover, they often ignore foundational testing principles, including boundary value analysis, equivalence partitioning, and exception handling, resulting in tests that may execute without error but fail to thoroughly validate functionality.

### 4.2 Knowledge-enhanced LLMs for SE Tasks

Recent studies have explored enhancing LLM performance on SE tasks by incorporating task-relevant context. In APR, Xia et al. [50] showed that simply providing the buggy function can outperform traditional tools, while ChatRepair [51] leveraged failing test names and assertions for interactive prompting. Further work has demonstrated the effectiveness of enriching context with bug-localized code [52], relevant identifiers [53, 54], stack traces [55], and bug reports [56].

In parallel, RAG has emerged as a promising paradigm for integrating external knowledge into LLMs by retrieving relevant information from codebases or databases. For code completion, ReACC [57] retrieves semantically similar code snippets, and Wu et al. [58] propose a selective RAG framework to reduce redundant retrievals. In code generation, recent approaches [59, 60] model code repositories as graph-based or knowledge-structured representations to retrieve and incorporate repository-level context into LLMs.

## 5 Threats to Validity

A key threat lies in the use of LLMs. To ensure fairness, all approaches are evaluated using the same model version. For baselines, we rely on official implementations or carefully reimplement them following published details, verifying outputs to align with reported results. Benchmark tasks and ground-truth tests are directly reused from prior work to ensure quality and consistency.

Another threat concerns generality and subjectivity. Our experiments focus on Java to remain consistent with prior work and enable direct comparison [20, 22], but KTester is not inherently tied to Java. Extending to other frameworks mainly requires updating prompts and replacing libraries, while supporting other languages involves only substituting the AST analysis tool (e.g., tree-sitter). Thus, the design is modular, and the adaptation cost is low. Future work will explore such extensions. While correctness and coverage provide objective metrics, they may not fully capture clarity and structure. We mitigate this via a user study on readability and maintainability, though results may still reflect participant bias. The study includes 14 professional developers, which we believe is sufficient, but broader sampling could further strengthen validity.

## 6 Conclusions

In this paper, we present KTester, a project-aware and testing-domain-knowledge-guided framework for LLM-based unit test generation. By extracting comprehensive project knowledge and decoupling test case design from test code generation, KTester effectively guides the LLM to produce semantically rich and maintainable test cases. Our multi-perspective prompting strategy ensures diverse and thorough coverage of testing scenarios, while structured test templates enhance code clarity and reusability. Extensive evaluations on real-world open-source projects demonstrate that KTester consistently outperforms existing state-of-the-art methods in correctness, coverage, readability, and maintainability. Future work includes exploring dynamic analysis to further enrich project knowledge and extending the approach to other programming languages and testing paradigms.

## References

[1] H. K. Leung and L. White, "Insights into regression testing (software testing)," in *Proceedings. Conference on Software Maintenance-1989*. IEEE, 1989, pp. 60–69.

[2] P. Runeson, "A survey of unit testing practices," *IEEE software*, vol. 23, no. 4, pp. 22–29, 2006.

[3] E. Daka and G. Fraser, "A survey on unit testing practices and problems," in *2014 IEEE 25th International Symposium on Software Reliability Engineering*. IEEE, 2014, pp. 201–211.

[4] E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo, "The oracle problem in software testing: A survey," *IEEE Trans. Software Eng.*, vol. 41, no. 5, pp. 507–525, 2015. [Online]. Available: https://doi.org/10.1109/TSE.2014.2372785

[5] G. Fraser and A. Arcuri, "A large-scale evaluation of automated unit test generation using evosuite," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 24, no. 2, pp. 1–42, 2014.

[6] M. Harman and P. McMinn, "A theoretical and empirical study of search-based testing: Local, global, and hybrid search," *IEEE Transactions on Software Engineering*, vol. 36, no. 2, pp. 226–247, 2009.

[7] A. Blasi, A. Gorla, M. D. Ernst, and M. Pezzè, "Call me maybe: Using nlp to automatically generate unit test cases respecting temporal constraints," in *37th IEEE/ACM International Conference on Automated Software Engineering*, 2022, pp. 1–11.

[8] P. Delgado-Pérez, A. Ramírez, K. J. Valle-Gómez, I. Medina-Bulo, and J. R. Romero, "Interevo-tr: Interactive evolutionary test generation with readability assessment," *IEEE Transactions on Software Engineering*, 2022.

[9] M. D. Ernst, J. H. Perkins, P. J. Guo, S. McCamant, C. Pacheco, M. S. Tschantz, and C. Xiao, "The daikon system for dynamic detection of likely invariants," *Science of computer programming*, vol. 69, no. 1-3, pp. 35–45, 2007.

[10] C. Csallner, N. Tillmann, and Y. Smaragdakis, "Dysy: Dynamic symbolic execution for invariant inference," in *Proceedings of the 30th international conference on Software engineering*, 2008, pp. 281–290.

[11] X. Xiao, S. Li, T. Xie, and N. Tillmann, "Characteristic studies of loop problems for structural test generation via symbolic execution," in *2013 28th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 2013, pp. 246–256.

[12] A. Zeller, R. Gopinath, M. Böhme, G. Fraser, and C. Holler, "The fuzzing book," 2019.

[13] C. Pacheco, S. K. Lahiri, M. D. Ernst, and T. Ball, "Feedback-directed random test generation," in *29th International Conference on Software Engineering (ICSE'07)*. IEEE, 2007, pp. 75–84.

[14] M. M. Almasi, H. Hemmati, G. Fraser, A. Arcuri, and J. Benefelds, "An industrial evaluation of unit test generation: Finding real faults in a financial application," in *2017 IEEE/ACM 39th International Conference on Software Engineering: Software Engineering in Practice Track (ICSE-SEIP)*. IEEE, 2017, pp. 263–272.

[15] "Ktester," 2025. [Online]. Available: https://github.com/SYSUSELab/KTester

[16] F. E. Allen, "Control flow analysis," *SIGPLAN Not.*, vol. 5, no. 7, p. 1–19, Jul. 1970. [Online]. Available: https://doi.org/10.1145/390013.808479

[17] "http://javaparser.org/," 2025.

[18] "https://spoon.gforge.inria.fr/about.html," 2025.

[19] 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*, ser. FSE 2024. New York, NY, USA: Association for Computing Machinery, Jul. 2024, pp. 572–576. [Online]. Available: https://dl.acm.org/doi/10.1145/3663529.3663801

[20] Z. Wang, K. Liu, G. Li, and Z. Jin, "HITS: High-coverage LLM-based Unit Test Generation via Method Slicing," in *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, ser. ASE '24. New York, NY, USA: Association for Computing Machinery, Oct. 2024, pp. 1258–1268. [Online]. Available: https://dl.acm.org/doi/10.1145/3691620.3695501

[21] P. Jaccard, "Lois de distribution florale dans la zone alpine," *Bulletin de la Société vaudoise des sciences naturelles*, vol. 38, pp. 69–130, 01 1902.

[22] Z. Yuan, Y. Lou, M. Liu, S. Ding, K. Wang, Y. Chen, and X. Peng, "Evaluating and improving chatgpt for unit test generation," *Proc. ACM Softw. Eng.*, vol. 1, no. FSE, 2024. [Online]. Available: https://doi.org/10.1145/3660783

[23] "https://github.com/zju-aces-ise/chatunitest-maven-plugin," 2025.

[24] "https://platform.openai.com," 2025.

[25] A. Deljouyi, R. Koohestani, M. Izadi, and A. Zaidman, *Leveraging Large Language Models for Enhancing the Understandability of Generated Unit Tests*. IEEE Press, 2025, p. 1449–1461. [Online]. Available: https://doi.org/10.1109/ICSE55347.2025.00032

[26] "Utgen replication package," 2025. [Online]. Available: https://github.com/amirdeljouyi/UTGen

[27] G. Fraser and A. Arcuri, "EvoSuite: automatic test suite generation for object-oriented software," in *Proceedings of the 19th ACM SIGSOFT symposium and the 13th European conference on Foundations of software engineering*, ser. ESEC/FSE '11. New York, NY, USA: Association for Computing Machinery, Sep. 2011, pp. 416–419. [Online]. Available: https://dl.acm.org/doi/10.1145/2025113.2025179

[28] "codellama-7b," 2025.

[29] S. Gu, Q. Zhang, K. Li, 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," Mar. 2025, arXiv:2408.03095 [cs]. [Online]. Available: http://arxiv.org/abs/2408.03095

[30] "https://www.jacoco.org/jacoco/," 2025.

[31] "claude-3-5-haiku-20241022 model overview," 2024. [Online]. Available: https://docs.claude.com/en/docs/about-claude/models/overview#legacy-models

[32] "deepseek-v3.1 release," 2025. [Online]. Available: https://api-docs.deepseek.com/news/news250821

[33] R. Pan, M. Kim, R. Krishna, R. Pavuluri, and S. Sinha, "ASTER: natural and multi-language unit test generation with llms," in *47th IEEE/ACM International Conference on Software Engineering: Software Engineering in Practice, SEIP@ICSE 2025, Ottawa, ON, Canada, April 27 - May 3, 2025*. IEEE, 2025, pp. 413–424. [Online]. Available: https://doi.org/10.1109/ICSE-SEIP66354.2025.00042

[34] M. Harman, Y. Jia, and Y. Zhang, "Achievements, open problems and challenges for search based software testing," in *2015 IEEE 8th international conference on software testing, verification and validation (ICST)*. IEEE, 2015, pp. 1–12.

[35] P. McMinn, "Search-based software testing: Past, present and future," in *2011 IEEE Fourth International Conference on Software Testing, Verification and Validation Workshops*. IEEE, 2011, pp. 153–163.

[36] R. Baldoni, E. Coppa, D. C. D'elia, C. Demetrescu, and I. Finocchi, "A survey of symbolic execution techniques," *ACM Computing Surveys (CSUR)*, vol. 51, no. 3, pp. 1–39, 2018.

[37] C. Cadar, D. Dunbar, D. R. Engler *et al.*, "Klee: unassisted and automatic generation of high-coverage tests for complex systems programs." in *OSDI*, vol. 8, 2008, pp. 209–224.

[38] S. K. Cha, T. Avgerinos, A. Rebert, and D. Brumley, "Unleashing mayhem on binary code," in *2012 IEEE Symposium on Security and Privacy*. IEEE, 2012, pp. 380–394.

[39] M. Tufano, D. Drain, A. Svyatkovskiy, S. K. Deng, and N. Sundaresan, "Unit test case generation with transformers and focal context," *arXiv preprint arXiv:2009.05617*, 2020.

[40] A. M. Dakhel, A. Nikanjam, V. Majdinasab, F. Khomh, and M. C. Desmarais, "Effective test generation using pre-trained large language models and mutation testing," *Information and Software Technology*, vol. 171, p. 107468, 2024.

[41] Y. Deng, C. S. Xia, C. Yang, S. D. Zhang, S. Yang, and L. Zhang, "Large language models are edge-case generators: Crafting unusual programs for fuzzing deep learning libraries," in *Proceedings of the 46th IEEE/ACM international conference on software engineering*, 2024, pp. 1–13.

[42] Y. Deng, C. S. Xia, H. Peng, C. Yang, and L. Zhang, "Large language models are zero-shot fuzzers: Fuzzing deep-learning libraries via large language models," in *Proceedings of the 32nd ACM SIGSOFT international symposium on software testing and analysis*, 2023, pp. 423–435.

[43] P. Nie, R. Banerjee, J. J. Li, R. J. Mooney, and M. Gligoric, "Learning deep semantics for test completion," *arXiv preprint arXiv:2302.10166*, 2023.

[44] J. Zhang, Y. Liu, P. Nie, J. J. Li, and M. Gligoric, "exlong: Generating exceptional behavior tests with large language models," *arXiv preprint arXiv:2405.14619*, 2024.

[45] 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*, vol. 50, no. 1, pp. 85–105, 2023.

[46] A. Sapozhnikov, M. Olsthoorn, A. Panichella, V. Kovalenko, and P. Derakhshanfar, "Testspark: Intellij idea's ultimate test generation companion," in *Proceedings of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion Proceedings*, 2024, pp. 30–34.

[47] C. Lemieux, J. P. Inala, S. K. Lahiri, and S. Sen, "Codamosa: Escaping coverage plateaus in test generation with pre-trained large language models," in *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2023, pp. 919–931.

[48] J. Altmayer Pizzorno and E. D. Berger, "Coverup: Effective high coverage test generation for python," *Proceedings of the ACM on Software Engineering*, vol. 2, no. FSE, pp. 2897–2919, 2025.

[49] G. Ryan, S. Jain, M. Shang, S. Wang, X. Ma, M. K. Ramanathan, and B. Ray, "Code-aware prompting: A study of coverage-guided test generation in regression setting using llm," *Proceedings of the ACM on Software Engineering*, vol. 1, no. FSE, pp. 951–971, 2024.

[50] C. S. Xia, Y. Wei, and L. Zhang, "Automated program repair in the era of large pre-trained language models," in *Proceedings of the 45th International Conference on Software Engineering*, ser. ICSE '23. IEEE Press, 2023, p. 1482–1494. [Online]. Available: https://doi.org/10.1109/ICSE48619.2023.00129

[51] C. S. Xia and L. Zhang, "Automated program repair via conversation: Fixing 162 out of 337 bugs for $0.42 each using chatgpt," in *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*, ser. ISSTA 2024. New York, NY, USA: Association for Computing Machinery, 2024, p. 819–831. [Online]. Available: https://doi.org/10.1145/3650212.3680323

[52] J. A. Prenner2024 and R. Robbes, "Out of context: How important is local context in neural program repair?" in *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, ser. ICSE '24. New York, NY, USA: Association for Computing Machinery, 2024. [Online]. Available: https://doi.org/10.1145/3597503.3639086

[53] C. S. Xia, Y. Ding, and L. Zhang, "The plastic surgery hypothesis in the era of large language models," in *Proceedings of the 38th IEEE/ACM International Conference on Automated Software Engineering*, ser. ASE '23. IEEE Press, 2024, p. 522–534. [Online]. Available: https://doi.org/10.1109/ASE56229.2023.00047

[54] N. Jiang, K. Liu, T. Lutellier, and L. Tan, "Impact of code language models on automated program repair," in *Proceedings of the 45th International Conference on Software Engineering*, ser. ICSE '23. IEEE Press, 2023, p. 1430–1442. [Online]. Available: https://doi.org/10.1109/ICSE48619.2023.00125

[55] M. Haque, P. Babkin, F. Farmahinifarahani, and M. Veloso, "Towards effectively leveraging execution traces for program repair with code LLMs," in *Proceedings of the 4th International Workshop on Knowledge-Augmented Methods for Natural Language Processing*, W. Shi, W. Yu, A. Asai, M. Jiang, G. Durrett, H. Hajishirzi, and L. Zettlemoyer, Eds. Albuquerque, New Mexico, USA: Association for Computational Linguistics, May 2025, pp. 160–179. [Online]. Available: https://aclanthology.org/2025.knowledgenlp-1.17/

[56] S. Fakhoury, S. Chakraborty, M. Musuvathi, and S. K. Lahiri, "Nl2fix: Generating functionally correct code edits from bug descriptions," in *Proceedings of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion Proceedings*, ser. ICSE-Companion '24. New York, NY, USA: Association for Computing Machinery, 2024, p. 410–411. [Online]. Available: https://doi.org/10.1145/3639478.3643526

[57] S. Lu, N. Duan, H. Han, D. Guo, S. won Hwang, and A. Svyatkovskiy, "Reacc: A retrieval-augmented code completion framework," 2022. [Online]. Available: https://arxiv.org/abs/2203.07722

[58] D. Wu, W. U. Ahmad, D. Zhang, M. K. Ramanathan, and X. Ma, "Repoformer: selective retrieval for repository-level code completion," in *Proceedings of the 41st International Conference on Machine Learning*, ser. ICML'24. JMLR.org, 2024.

[59] J. Li, X. Shi, K. Zhang, L. Li, G. Li, Z. Tao, J. Li, F. Liu, C. Tao, and Z. Jin, "Coderag: Supportive code retrieval on bigraph for real-world code generation," 2025. [Online]. Available: https://arxiv.org/abs/2504.10046

[60] M. Athale and V. Vaddina, "Knowledge graph based repository-level code generation," in *2025 IEEE/ACM International Workshop on Large Language Models for Code (LLM4Code)*. IEEE, May 2025, p. 169–176. [Online]. Available: http://dx.doi.org/10.1109/LLM4Code66737.2025.00026