# Enhancing LLM's Ability to Generate More Repository-Aware Unit Tests Through Precise Contextual Information Injection

XIN YIN, The State Key Laboratory of Blockchain and Data Security, Zhejiang University, China
CHAO NI*, The State Key Laboratory of Blockchain and Data Security, Zhejiang University, China
XINRUI LI, The State Key Laboratory of Blockchain and Data Security, Zhejiang University, China
LIUSHAN CHEN, ByteDance Inc., China
GUOJUN MA, ByteDance Inc., China
XIAOHU YANG, The State Key Laboratory of Blockchain and Data Security, Zhejiang University, China

Though many learning-based approaches have been proposed for unit test generation and achieved remarkable performance, they still have limitations in relying on task-specific datasets. Recently, Large Language Models (LLMs) guided by prompt engineering have gained attention for their ability to handle a broad range of tasks, including unit test generation. Despite their success, LLMs may exhibit hallucinations when generating unit tests for focal methods or functions due to their lack of awareness regarding the project's global context. These hallucinations may manifest as calls to non-existent methods, as well as incorrect parameters or return values, such as mismatched parameter types or numbers. While many studies have explored the role of context, they often extract fixed patterns of context for different models and focal methods, which may not be suitable for all generation processes (e.g., excessive irrelevant context could lead to redundancy, preventing the model from focusing on essential information).

To overcome this limitation, we propose RATester, which enhances the LLM's ability to generate more repository-aware unit tests through global contextual information injection. To equip LLMs with global knowledge similar to that of human testers, we integrate the language server gopls, which provides essential features (e.g., definition lookup) to assist the LLM. When RATester encounters an unfamiliar identifier (e.g., an unfamiliar struct name), it first leverages gopls to fetch relevant definitions and documentation comments, and then uses this global knowledge to guide the LLM. By utilizing gopls, RATester enriches the LLM's knowledge of the project's global context, thereby reducing hallucinations during unit test generation.

We evaluate the effectiveness and efficiency of RATester compared to baseline approaches by constructing a new Golang dataset from real-world projects. The results demonstrate the advantages of RATester over the baselines. For instance, on our dataset, RATester achieves an average line coverage of 26.25%, representing an improvement of 16.30% to 165.69% over the baselines. Furthermore, RATester shows superior performance in mutation testing, successfully killing 25 to 147 more mutants than the baseline approaches. Additionally, we extend our analysis to assess the model-agnostic effectiveness of RATester. These results not only confirm the effectiveness of RATester but also underscore its universal applicability.

*Chao Ni the corresponding author.
He is also with Hangzhou High-Tech Zone (Binjiang) Institute of Blockchain and Data Security.

Authors' Contact Information: Xin Yin, The State Key Laboratory of Blockchain and Data Security, Zhejiang University, Hangzhou, China, xyin@zju.edu.cn; Chao Ni, The State Key Laboratory of Blockchain and Data Security, Zhejiang University, Hangzhou, China, chaoni@zju.edu.cn; Xinrui Li, The State Key Laboratory of Blockchain and Data Security, Zhejiang University, Hangzhou, China, lixinrui@zju.edu.cn; Liushan Chen, ByteDance Inc., Shenzhen, China, chenliushan@bytedance.com; Guojun Ma, ByteDance Inc., Shenzhen, China, maguojun@bytedance.com; Xiaohu Yang, The State Key Laboratory of Blockchain and Data Security, Zhejiang University, Hangzhou, China, yangxh@zju.edu.cn.

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

Additional Key Words and Phrases: Unit Test Generation, Large Language Model, Global Context

## 1 Introduction

Unit testing plays a critical role in software maintenance by enabling developers to identify defects and errors early in the development process, thereby ensuring the quality of software systems. This not only helps lower overall product costs but also enhances developer productivity [8, 16, 20]. Despite its significance, manually writing high-quality unit tests is both challenging and time-consuming. To mitigate this, a range of automated unit test generation approaches have been proposed, which can be broadly classified into three categories: traditional approaches [13, 29, 44], learning-based approaches [6, 18, 34, 40, 42], and LLM-based approaches [9, 21, 24, 54].

Traditional approaches [13, 29] primarily focus on maximizing code coverage, and research demonstrates their effectiveness in achieving high coverage [7, 26, 30, 31]. Randoop [29] and EvoSuite [13] are among the most popular and widely used examples of such approaches. Randoop, a widely recognized tool, is extensively used to generate unit tests for Java code through feedback-directed random test generation. EvoSuite automates the creation of test suites, aiming to maximize code coverage while minimizing test suite size and ensuring comprehensive assertions. Despite their success in achieving coverage goals, previous studies show that these approaches do not produce well-written, maintainable unit tests explicitly for developers to use [42, 44].

To overcome this limitation, learning-based approaches [6, 18, 34, 40, 42] have gained significant attention in recent years. AthenaTest [42] utilizes a Transformer-based model trained on developer-written test cases to generate accurate and readable unit tests. A3Test [6] applies domain adaptation techniques, aiming to transfer knowledge from assertion generation tasks to test case generation. UniTester [18] leverages the UniTSyn dataset to synthesize unit tests across multiple programming languages, including Golang. However, these approaches treat unit test generation as a translation problem, where the goal is to translate a focal method or function into a corresponding unit test. They rely heavily on task-specific datasets extracted from open-source repositories.

To address the challenges faced by learning-based approaches, researchers are increasingly exploring pre-trained Large Language Models (LLMs) for unit test generation. These models generate unit tests directly from contextual information, reducing reliance on task-specific datasets by leveraging extensive pre-training on large open-source code snippets. Researchers [24, 54] have adopted ChatGPT to generate unit tests based on focal methods. Despite these advancements, LLMs can still exhibit hallucinations when generating unit tests for focal methods or functions due to their lack of awareness regarding the project's global context. These hallucinations can include, but are not limited to, calling non-existent methods, as well as assigning incorrect parameters and return values (e.g., mismatched parameter types or incorrect parameter numbers). To overcome this limitation, many studies have explored the extraction of context to reduce hallucinations in the generation process of LLMs. ChatUniTest [9] introduces an LLM-based framework that enhances automated unit test generation with an adaptive focal context mechanism, capturing relevant context within prompts. It also employs a "Generation-Validation-Repair" process to correct errors in the generated tests. Following that, researchers [15, 36, 53] have explored the roles of focal context and dependency context. These methods utilize one or more fixed patterns to extract context for the focal method: (1) focal class signature; (2) signatures of other methods and fields

in the class; (3) signatures of dependent classes; and (4) signatures of dependent methods and fields in the dependent classes. However, these fixed extraction patterns present several issues: (1) they may overlook important context; for instance, when generating a unit test for a specific focal method, the LLM might require unknown context beyond the dependencies of that focal method; (2) there is potential for redundant context, as excessive irrelevant context could lead to redundancy, preventing the model from focusing on essential information.

In practical development scenarios, developers are typically highly familiar with the methods, functions, and structs within the package they are working on. Additionally, Integrated Development Environment (IDE) tools and language servers further assist by providing information on function calls and identifier descriptions, enabling developers to produce more accurate code. Therefore, in this paper, we aim to provide LLMs with a project's global knowledge comparable to that of human testers by introducing RATester, which enhances LLM's ability to generate more repository-aware unit tests through global contextual information injection. **First, we propose an LLM-based framework for unit test generation.** LLMs are trained in an unsupervised manner using up to billions of text and code tokens. This extensive unsupervised learning process equips LLMs with robust reasoning capabilities, enabling them to generate unit tests without relying on task-specific training datasets. Therefore, we propose a novel LLM-based approach RATester for unit test generation since the representative conversational LLM provides advanced capabilities for several tasks, including natural language processing [27], code generation [22], and unit test generation [9, 24, 54]. **Second, we develop a global-aware framework to enhance the capabilities of LLMs.** To provide LLMs with a global knowledge similar to that of human testers, we integrate the language server gopls [1], which provides features such as code completion, syntax checking, and definition lookup, significantly improving development efficiency. When RATester encounters unfamiliar identifiers (e.g., unfamiliar method names, unfamiliar function names, and unfamiliar struct names), it proactively invokes gopls to fetch relevant definitions and documentation comments. By continuously leveraging the capabilities of gopls, RATester progressively enriches the LLM's global knowledge of the project, thereby reducing hallucinations and improving the effectiveness of unit test generation.

We construct a dataset to evaluate RATester, consisting of eight highly starred GitHub projects (with stars ranging from 29.7k to 85.5k): beego, echo, fiber, frp, gin, hugo, nps, and traefik. To evaluate the effectiveness and efficiency of RATester, we compare it against five baseline approaches across three categories: one traditional approach (i.e., NxtUnit [44]), one learning-based approach (i.e., UniTester [18]), and three basic LLMs (i.e., CodeLlama [35], DeepSeek-Coder [5], and Magicoder [46]). The results demonstrate the clear superiority of RATester over the baselines. For instance, on our collected dataset, RATester achieves an average line coverage of 26.25%, representing an improvement of 16.30% to 165.69% over the baselines. Furthermore, RATester achieves the highest performance in mutation testing, successfully killing 25 to 147 more mutants than the baseline approaches. We also extend our analysis to explore the model-agnostic capabilities of RATester. The results not only validate the effectiveness of RATester but also emphasize its universal applicability. RATester is designed to be model-agnostic, enabling it to adapt to various LLMs, further underscoring its flexibility and universality.

In summary, the key contributions of this paper include:

**A. Novel LLM-based Framework:** We present RATester, an advanced LLM-based framework for unit test generation that does not rely on task-specific training datasets. Our results demonstrate that this framework can outperform existing approaches, achieving superior performance in unit test generation.

**B. Repository-Aware Tester:** We introduce RATester, which utilizes the features (e.g., definition lookup) of gopls to enhance the LLM's global knowledge of the project. By proactively

fetching definitions and documentation comments for unfamiliar information, RATester reduces hallucinations during unit test generation.

**C. Extensive Evaluation:** (1) We conduct studies on the effectiveness and efficiency of RATester and baselines by collecting a new Golang dataset from real-world projects. (2) We evaluate RATester and baselines not only using compile rate and line coverage metrics but also assess their capabilities in mutation testing.

## 2   Motivation

### 2.1   A Motivation Example

Fig. 1 shows a focal method named "PATCH" along with the unit tests generated by DeepSeek-Coder for a Golang project named gin. The upper right corner of Fig. 1 illustrates how DeepSeek-Coder (using imprecise context) generates a unit test for the focal method without global knowledge of the project. The unit test "TestPATCH" verifies whether the server can correctly handle an HTTP PATCH request sent to the "/patch" path and return the expected response status code of "http.StatusOK" and the response body "Hello, World". The test creates a route instance, defines a handler for the PATCH request, and then uses the "httptest" package to simulate the request and capture the response, ultimately checking whether the response's status code and body meet expectations. However, in the fourth line, this unit test encounters a compilation error: "c.String(http.StatusOK, 'Hello, World') (no value) used as value", preventing the test from compiling. This issue arises because DeepSeek-Coder (using imprecise context) lacks sufficient knowledge of the project and does not know that the "String" method within the "Context" struct does not return a value, leading to hallucinations during inference.

To generate a correct unit test, one needs to understand the definition and documentation comments for the "String" method in the "Context" struct to call it correctly. In real-world development scenarios, human testers are typically well-acquainted with the project's global context, including the methods, functions, and structs within the package. IDE tools or language servers further enhance this familiarity by providing call information and identifier descriptions, aiding human testers in writing accurate code. As a result, human testers often refer to the tips provided by these tools to supplement their global knowledge when crafting unit tests for focal methods or functions, thereby reducing the occurrence of erroneous test cases.

**Observation.** Due to the use of fixed patterns for context extraction in existing methods, there is a certain degree of knowledge omission and redundancy in the context provided to LLMs, meaning it may not be the context that LLMs truly need during the generation process. Given the limited input size of LLMs, it is impossible to feed all information into the model. Excessive irrelevant context can lead to redundancy, preventing the model from focusing on essential information, which may result in LLMs exhibiting hallucinations when generating unit tests for focal methods or functions. These hallucinations may include invoking non-existent methods, setting incorrect parameters and return values (e.g., parameter type mismatches or incorrect number of parameters). We should provide the model with efficient information necessary for generation, rather than fixed selections, to reduce the interference of redundant information.

In real-world development scenarios, human testers are typically very familiar with the package, which helps reduce the occurrence of such hallucinations. Therefore, we also conduct a foundational experiment to simulate human testers' global knowledge of the package by providing the definition and documentation comments for the "Context" struct as input above the focal method. In this case, we find that DeepSeek-Coder (using precise context) successfully generates a compilable unit test, as shown on the lower right corner of Fig. 1. In the fifth line, DeepSeek-Coder (using precise context) correctly uses the "String" method within the "Context" struct without attempting to set a

return value. We believe this results from DeepSeek-Coder's foundational understanding of the "Context" struct, which prevents it from generating erroneous usages. This foundational knowledge reflects DeepSeek-Coder's understanding of the project, similar to a human tester's familiarity with it. High-quality foundational knowledge enables the model to generate higher-quality unit tests.
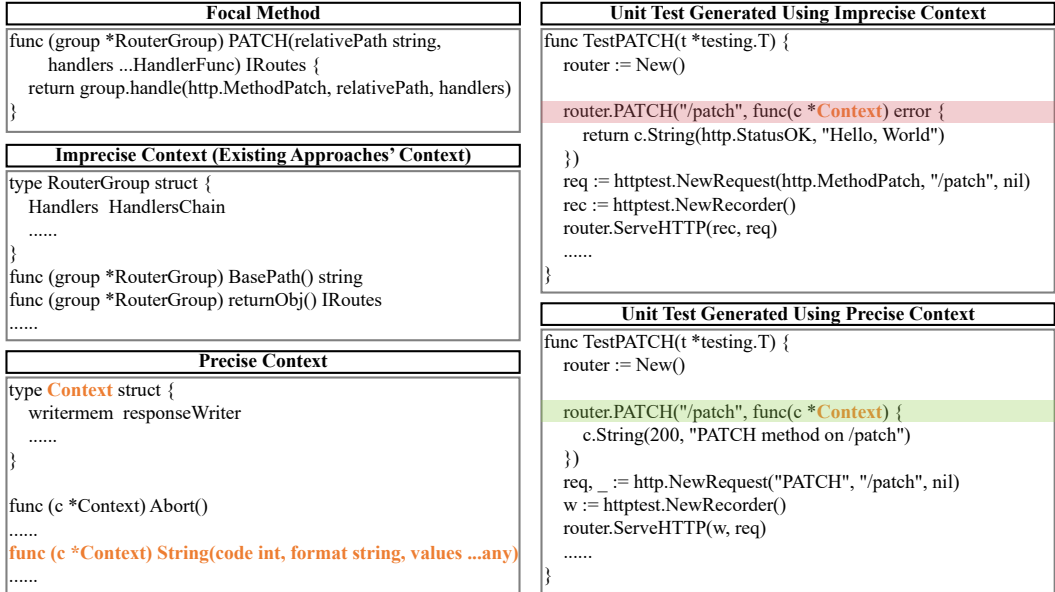
| Focal Method |
|---|
| func (group *RouterGroup) PATCH(relativePath string,<br>    handlers ...HandlerFunc) IRoutes {<br>    return group.handle(http.MethodPatch, relativePath, handlers)<br>} |

| Imprecise Context (Existing Approaches' Context) |
|---|
| type RouterGroup struct {<br>    Handlers  HandlersChain<br>    ......<br>}<br>func (group *RouterGroup) BasePath() string<br>func (group *RouterGroup) returnObj() IRoutes<br>...... |

| Precise Context |
|---|
| type **Context** struct {<br>    writermem  responseWriter<br>    ......<br>}<br><br>func (c *Context) Abort()<br>......<br>**func (c *Context) String(code int, format string, values ...any)**<br>...... |

| Unit Test Generated Using Imprecise Context |
|---|
| func TestPATCH(t *testing.T) {<br>    router := New()<br><br>    router.PATCH("/patch", func(c ***Context**) error {<br>        return c.String(http.StatusOK, "Hello, World")<br>    })<br>    req := httptest.NewRequest(http.MethodPatch, "/patch", nil)<br>    rec := httptest.NewRecorder()<br>    router.ServeHTTP(rec, req)<br>    ......<br>} |

| Unit Test Generated Using Precise Context |
|---|
| func TestPATCH(t *testing.T) {<br>    router := New()<br><br>    router.PATCH("/patch", func(c ***Context**) {<br>        c.String(200, "PATCH method on /patch")<br>    })<br>    req, _ := http.NewRequest("PATCH", "/patch", nil)<br>    w := httptest.NewRecorder()<br>    router.ServeHTTP(w, req)<br>    ......<br>} |

**Fig. 1.** A focal method along with the unit tests generated by DeepSeek-Coder for a project named gin

## 2.2 Key Ideas

Based on the above observation, we propose RATester for unit test generation that utilizes gopls to fetch definitions and documentation comments for unfamiliar structs, methods, functions, and more. This approach provides LLMs with the precise context needed for generation, effectively reducing hallucinations during unit test generation.

(1) **LLM-based generation approach.** Unlike learning-based approaches (e.g., TOGA [11], A3Test [6], and UniTester [18]), LLMs are trained in an unsupervised manner using up to billions of text and code tokens. This large-scale unsupervised learning process allows LLMs to have strong reasoning capabilities and be applied for unit test generation without relying on training with a large amount of task-specific data. Therefore, we propose an LLM-based generation approach, namely RATester, since the representative conversational LLM provides advanced capabilities for several tasks [47–51], including unit test generation [24, 38].

(2) **Proactively fetch global contextual information.** Due to the limited input size of LLMs, it is impossible to feed all information into them. Consequently, existing approaches use fixed patterns to select context information to include in the prompts. However, these fixed selection approaches suffer from issues of knowledge omission and redundancy. In this paper, to endow LLMs with a global knowledge similar to that of human testers, we introduce the Golang language server gopls, which can provide features (e.g., definition lookup) for LLMs. For example, as illustrated on the right side of Fig. 1, when RATester encounters the unfamiliar method "Context", it proactively uses gopls to fetch specific definitions and documentation comments to avoid erroneous usage. Therefore, we propose RATester, an global-aware tester for unit test generation that utilizes gopls

to fetch definitions and documentation comments for unfamiliar structs, methods, functions, and more. By leveraging the capabilities of gopls, we aim to enhance the LLM's global understanding of the project and reduce hallucinations during unit test generation.

## 3 Our Approach: RATester

In this section, we present the methodology behind our RATester approach. We begin with an overview of the approach, followed by a detailed discussion of each component.

### 3.1 Overview

As shown in Fig. 2, our approach consists of three components: Fetcher, Formulator, and Generator. Given the focal method or function that needs to be tested (referred to as the method or function under test), each component plays a distinct role in the unit test generation task:

- **Fetcher** fetches the definition and documentation comment of unfamiliar identifier (e.g., struct name) using language server (e.g., gopls) according to the given information.
- **Formulator** fills the fetched code context and the test snippet with newly generated identifiers into the prompt template, and then formulates them as input for the generator.
- **Generator** leverages the formulated input to perform the unit test generation. We use DeepSeek-Coder [5] as the generator, which can be replaced with various LLMs (e.g., CodeLlama [35]).
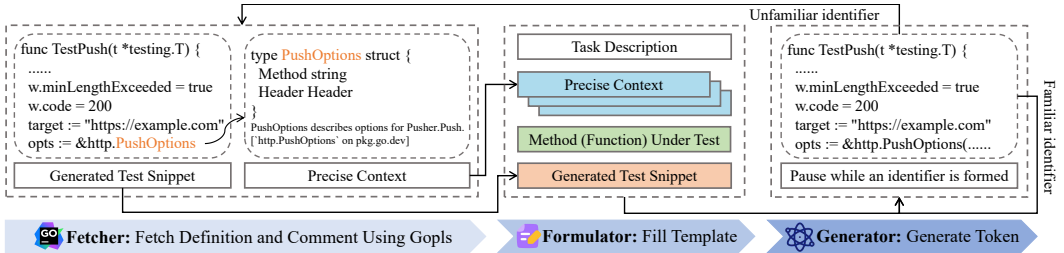


**Fig. 2.** Overview of RATester

### 3.2 Fetcher

When generating unit tests for focal methods or functions, LLMs often produce hallucinations, which can manifest as calls to non-existent methods, as well as incorrect parameter assignments and return values, such as mismatched parameter types or an improper number of parameters. In contrast, human testers typically possess a strong understanding of the various methods, functions, and structs within the package during the test development process. Additionally, IDE tools and language servers provide essential support by offering information on function calls and identifier descriptions, facilitating the creation of accurate code. Consequently, human testers frequently leverage insights from these tools to enhance their global knowledge while crafting unit tests, ultimately reducing the likelihood of erroneous test cases.

In this paper, RATester serves as a Fetcher by utilizing gopls to provide LLMs with a global knowledge comparable to that of human testers. Gopls [1], the Go language server, facilitates interactions with editors such as Visual Studio Code. By leveraging the capabilities of gopls, we can improve the LLM's global understanding of the package by providing precise context, thereby mitigating hallucinations during the generation of unit tests.

In the initial stage of unit test generation, RATester actively queries gopls for the definitions and documentation comments of the receiver type (which does not exist if a function is being tested), the parameter types, and the return type of the focal method or function. All fetched information is

then filled into the prompt template. During the continuous phase of unit test generation, whenever the LLM generates an unfamiliar identifier (e.g., new function name and new struct name), RATester proactively utilizes gopls to check whether the identifier exists in the current package and fetches its definition and documentation comments. After this, the fetched information is also filled into the prompt template by the Formulator (refer to Section 3.3 for more details). By leveraging gopls to proactively fetch definitions and documentation comments at both stages and enriching the prompt with this knowledge, the LLM gains a comprehensive understanding of unfamiliar information. This process closely resembles how human testers utilize IDE tools to look up definitions and documentation for unfamiliar methods, functions, and structs. Our RATester approach enables the LLM to act as a tester with extensive project knowledge, thereby enhancing its testing capabilities.
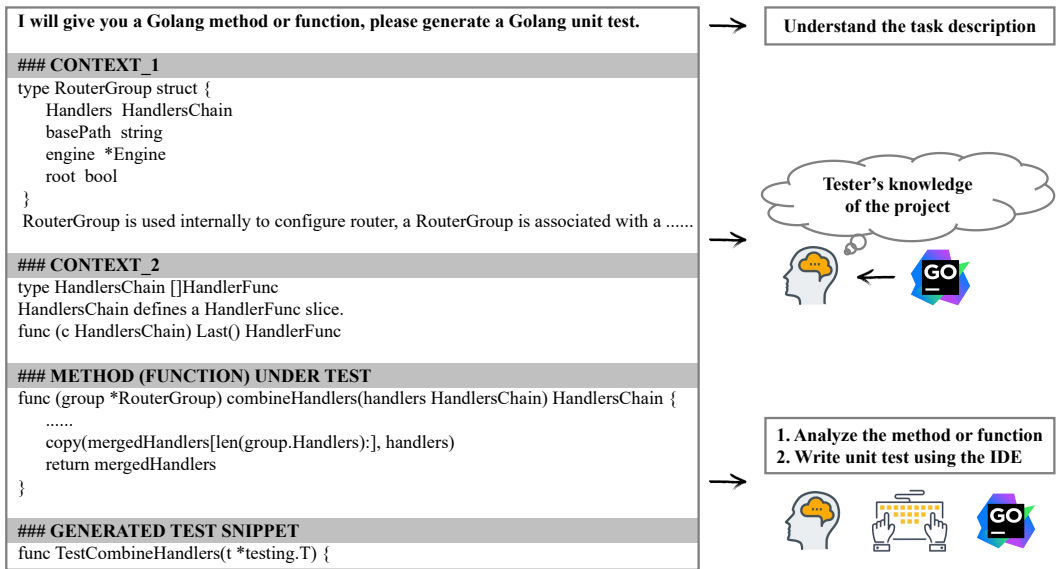


**Fig. 3.** An example of prompt for the unit test generation

### 3.3 Formulator

In both the initial stage and the continuous phase of unit test generation, the formulator fills the fetched code context and the test snippet into the prompt template, where fill means the procedure of inserting content into corresponding positions of the template. As shown in Fig. 3, this prompt template consists of four main parts:

- **Task Description.** RATester provides the LLM with the description constructed as *"I will give you a Golang method or function, please generate a Golang unit test"*. This part aids the LLM in understanding the task description, simulating the process by which a human tester comprehends the objectives of the task.
- **Precise Context.** RATester provides the fetched definitions and documentation comments to LLM. The code context continuously expands as the generation process progresses, enhancing the LLM's global knowledge of the project. This part simulates the global project knowledge that human testers possess with the assistance of IDE tools.
- **Method (Function) Under Test.** RATester provides the focal method or function to LLM. We also prefix the focal method or function with *"### METHOD (FUNCTION) UNDER TEST"* to

directly indicate LLM about the context of the method or function. This part simulates the scenario in which human testers review the method or function being tested.

- **Generated Test Snippet**. RATester provides the LLM with the unit test generated in the previous round, along with a new identifier. In the initial stage, the generated test snippet is explicitly set as *"func Test{name}(t *testing.T)"*. As this part continues to expand, it simulates the iterative process of human testers writing unit tests.

---

**Algorithm 1:** Unit Test Generation Process

---

**Input:** Focal method (function) $M$, Focal method (function) name $N$, Prompt template $P$;

**Initialize:** list_i = [], token_length = 0, test_snippet = "func Test{$N$}(t *testing.T) {";

*1. Fetch definitions and documentation comments of M (Section 3.2);*

*2. Append definitions and documentation comments to code context;*

*3. Fill code context and test snippet into P (Section 3.3);*

**while** *token_length++ < 512* **do**
  *Generate token t using P;*
  *test_snippet += t;*
  **if** *is_golang_identifier_part(t)* **then**
    │ *Append t to list_i;*
  **else**
    │ **if** *list_i is not empty* **then**
    │   │ *identifier = concatenate(list_i);*
    │   │ **if** *identifier not in code context* **then**
    │   │   │ *1. Fetch definitions and documentation comments of identifier;*
    │   │   │ *2. Append definitions and documentation comments to code context;*
    │   │ *Set list_i to empty;*
  *Fill code context and test snippet into P;*

**Output:** All generated tokens;

---

## 3.4 Generator

The generator leverages results returned from the formulator and performs the tasks of unit test generation accordingly. It continually generates tokens until a complete identifier is formed. For unfamiliar identifiers, RATester actively invokes the Fetcher to supplement the code context, while for familiar identifiers, RATester fills the new identifier into the prompt to generate next token. As shown in Algorithm 1, RATester initializes a List: list_i to store the currently generated identifiers. If the token being generated can be part of an identifier (i.e., it follows Golang's identifier naming rules), it is added to list_i. If the token cannot be part of an identifier (e.g., the LLM generates a character like '.'), and list_i is not empty, the process considers that a complete identifier has been generated. If the generated identifier is not found in the code context, RATester uses gopls to fetch the identifier's definition and documentation comments, which are then filled into the prompt template for the next generation step (refer to Section 3.2 and Section 3.3 for more details). By continuously leveraging gopls to fetch the code context, the LLM acquires sufficient global knowledge, thereby reducing the likelihood of hallucinations during the generation process. In this paper, we adopt DeepSeek-Coder [5] as the backend LLM. RATester is flexible to include other LLMs as the backend model (e.g., CodeLlama [35] and Magicoder [46]).

## 4 Experimental Design

In this section, we first present our collected dataset and then introduce the baseline approaches. Following that, we describe the performance metrics as well as the experimental setting.

**Table 1.** The statistic of constructed dataset

| Project | Star | Focal Method and Function (#) | Line Coverage of Unit Tests (%) |
|---|---|---|---|
| **beego** | 31.5k | 2,688 | 38.78% |
| **echo** | 29.7k | 419 | 93.58% |
| **fiber** | 33.6k | 765 | 85.46% |
| **frp** | 85.5k | 864 | 2.59% |
| **gin** | 78.5k | 449 | 95.53% |
| **hugo** | 75.4k | 3,829 | 76.54% |
| **nps** | 30.6k | 455 | 0.51% |
| **traefik** | 50.9k | 1,726 | 58.95% |

### 4.1 Dataset Construction

We construct a dataset to evaluate RATester, consisting of eight highly-starred GitHub projects (with stars ranging from 29.7k to 85.5k): beego, echo, fiber, frp, gin, hugo, nps, and traefik. Since RATester focuses on generating unit tests for focal methods and functions, we extract all methods and functions from each project. The detailed information for each project is displayed in Table 1. In addition to the star count and the number of successfully extracted focal methods and functions, we also run all unit tests within the projects and show the line coverage.

### 4.2 Baselines

To investigate the effectiveness of RATester, we consider six baselines for a comprehensive comparison, including one traditional approach, one learning-based approach, and four LLM-based approaches:

**Traditional Approach.** To present the traditional approach, we employ NxtUnit [44], an automatic unit test generation tool for Go that leverages random testing and is particularly suited for microservice architectures. It offers three types of interfaces: an integrated development environment (IDE) plugin, a command-line interface (CLI), and a web-based platform. NxtUnit's random-based strategy allows it to quickly generate unit tests, making it ideal for smoke testing and rapid quality feedback. However, NxtUnit may sometimes fail to generate test cases due to issues like compilation errors or test crashes. As a result, NxtUnit only provides test cases that can be executed successfully.

**Learning-based Approach.** To present the learning-based approach, we utilize the transformer-based generation model, UniTester [18]. This model is trained on the UniTSyn dataset and is capable of synthesizing unit tests for programs in multiple languages, including Golang. As the published code for UniTester lacks the model checkpoint, we retrain the UniTester model following the settings described in the paper and using the UniTSyn dataset. To prevent data leakage, we exclude projects from the training set that overlap with those in our collected dataset.

**LLM-based Approach.** To represent the LLM-based approach, we utilize ChatUniTest [9] and three LLMs to generate unit tests for each focal method and function without fine-tuning. The models we select are recently released: CodeLlama [35], DeepSeek-Coder [5], and Magicoder [46].

- **CodeLlama** proposed by Rozière et al. [35] is a collection of large pre-trained language models for code, built on Llama 2 architecture. These models achieve state-of-the-art among open-source

models for code-related tasks, offer infilling capabilities, support for large input contexts, and robust zero-shot instruction-following abilities for programming problems.

- **DeepSeek-Coder** developed by DeepSeek AI [5] comprises a series of code language models, each trained from scratch on 2 trillion tokens. The training data includes 87% code and 13% natural language, covering both English and Chinese. DeepSeek-Coder demonstrates state-of-the-art performance among open-source code models, excelling across multiple programming languages and various benchmarks.
- **Magicoder** proposed by Wei et al. [46] is trained on 75K synthetic instruction data using OSS-Instruct, a novel approach that leverages open-source code snippets to generate diverse instruction data for code-related tasks. The approach aims to mitigate the inherent bias in LLM-generated synthetic data by harnessing the vast resources of open-source code, leading to more realistic and controllable data generation.

### 4.3 Evaluation Metrics

To evaluate the performance of RATester and baseline approaches, we use Compile Rate and Line Coverage as primary metrics:

**Compile Rate** represents the proportion of test cases that can be successfully compiled and executed out of the total number generated. A higher compile rate reflects better quality and reliability in the generated test cases.

**Line Coverage** quantifies the percentage of code lines executed by the test cases, offering insights into the effectiveness of the tests in covering different parts of the code. A higher line coverage indicates that a larger portion of the code is being tested.

While Compile Rate and Line Coverage are valuable, they do not fully assess test quality. To provide a more comprehensive evaluation of the unit tests generated by RATester, we also employ mutation testing. We use Gremlins [2] to introduce mutations into the projects and evaluate the number of mutants killed by the unit tests, along with mutator coverage.

### 4.4 Implementation Details

We develop the unit test generation in Python, utilizing PyTorch [32] implementations of LLMs (i.e., CodeLlama 7B, DeepSeek-Coder 6.7B, and Magicoder 6.7B). We use the Hugging Face API [3] to load the model weights and generate outputs. We also adhere to the best-practice guide [39] for our prompt. For the baseline comparisons, we directly use the settings provided in the NxtUnit's [44] original paper to generate unit tests. Since the published code for UniTester [18] does not include the model checkpoint, we retrain the UniTester model using the settings and UniTSyn dataset provided in the original paper. To avoid data leakage, we exclude any data from the UniTSyn dataset that overlaps with those in our collected dataset during training. Considering both the performance improvements and the associated generation costs, we generate one unit test for each focal method and function (refer to Section 5.3 for more details) and test them using the go test command. Our evaluation is conducted on a 32-core workstation equipped with an Intel(R) Xeon(R) Platinum 8358P CPU @ 2.60GHz, 2TB RAM, and 8×NVIDIA A800 80GB GPU, running Ubuntu 20.04.6 LTS.

## 5 Experimental Results

To investigate the effectiveness of RATester on unit test generation, our experiments focus on the following three research questions:

- **RQ-1 Effectiveness Comparison.** *How does the performance of RATester compare with the baselines in unit test generation?*

- **RQ-2 Model-Agnostic Analysis.** *What are the model-agnostic capabilities of RATester in unit test generation?*
- **RQ-3 Efficiency Comparison.** *How does the efficiency of RATester compare with the baselines in unit test generation?*

### 5.1 RQ-1: Effectiveness of RATester

**Objective.** To reduce the hallucination issues that LLMs experience during unit test generation (e.g., invoking non-existent methods and setting incorrect parameters and return values), we propose the RATester approach. This approach utilizes the definition lookup feature provided by gopls to dynamically fetch relevant code context during the generation process. By supplying LLMs with more project-specific knowledge, we aim to minimize hallucinations. In this section, our objective is to investigate whether RATester outperforms previous unit test generation approaches in terms of effectiveness.

**Experimental Design.** In this RQ, we employ DeepSeek-Coder as the backend model for RATester. To facilitate a fair comparison, we consider three baselines: NxtUnit [44], UniTester [18], and ChatUniTest [9]. For NxtUnit, we use the default settings from the original paper. For ChatUniTest, we utilize DeepSeek-Coder as the backend model. As the published code for UniTester lacks the model checkpoint, we retrain the UniTester model using the settings and UniTSyn dataset from the original paper. To prevent data leakage, we exclude any overlapping data between the UniTSyn and our collected datasets during training.

For a comprehensive performance comparison between the baselines and RATester, we conduct two distinct experiments across eight real-world projects. The first experiment involves executing all generated unit tests within each project, recording the compile rate and line coverage. In the second experiment, we extend our evaluation with mutation testing, using Gremlins [2] to mutate the projects and measure the number of mutants killed by the generated unit tests, along with the mutator coverage. This experiment demonstrates the effectiveness of unit tests generated by RATester in detecting unknown defects.

**Results.** We discuss the results from the aspects of compile rate, line coverage, and mutation testing, respectively.

**Effectiveness of RATester in Compile Rate and Line Coverage.** Table 2 shows the compile rate and line coverage of unit tests generated by RATester and the baselines. We observe that RATester consistently outperforms the baselines across all projects. Specifically, the compile rate of unit tests generated by RATester significantly improves from 16.67%−63.56% to 45.58%−69.49%, with the average compile rate increasing from 24.61% to 61.84%. Note that NxtUnit only provides test cases that can be executed successfully. Therefore, we do not record the compile rate for the test cases generated by NxtUnit.

In addition to the compile rate, RATester demonstrates a significant enhancement in line coverage. Across all evaluated projects, RATester increases the line coverage from a range of 7.49%−53.92% to 12.92%−58.09%. This results in an average improvement of 16.30%−165.69% when compared to the baseline approaches. Such a substantial increase in line coverage indicates that RATester is more effective in generating comprehensive unit tests, thereby enhancing the overall robustness of the tested software.

In Table 3, we evaluate the ability of unit tests generated by RATester and the baselines to enhance the line coverage of the original unit tests across the eight projects. The "Original" column represents the line coverage of the original unit tests. In contrast, the "Original+NxtUnit" column displays the total coverage achieved by combining the original tests with those generated by NxtUnit. Similarly, the "Original+UniTester" column illustrates the total coverage obtained by integrating the

**Table 2.** RQ-1: RATester vs. Baselines across different projects in compile rate and line coverage

| Projects | Compile Rate | | | | Line Coverage | | | |
|---|---|---|---|---|---|---|---|---|
| | NxtUnit | UniTester | ChatUniTest | RATester | NxtUnit | UniTester | ChatUniTest | RATester |
| **beego** | - | 31.87% | 56.38% | 68.75% | 20.51% | 12.71% | 27.85% | 31.91% |
| **echo** | - | 22.36% | 39.65% | 45.58% | 10.77% | 9.89% | 23.09% | 24.50% |
| **fiber** | - | 21.22% | 42.44% | 59.61% | 20.24% | 9.33% | 19.33% | 26.31% |
| **frp** | - | 28.57% | 52.35% | 66.90% | 12.84% | 8.62% | 13.20% | 14.43% |
| **gin** | - | 32.91% | 63.56% | 69.49% | 21.92% | 11.38% | 53.92% | 58.09% |
| **hugo** | - | 24.87% | 48.76% | 61.51% | 12.34% | 9.57% | 18.78% | 25.02% |
| **nps** | - | 16.67% | 53.11% | 64.84% | 16.48% | 7.49% | 12.66% | 16.82% |
| **traefik** | - | 18.44% | 46.25% | 58.05% | 10.45% | 10.01% | 11.69% | 12.92% |
| **Average** | - | 24.61% | 50.31% | 61.84% | 15.69% | 9.88% | 22.57% | 26.25% |

original tests with those produced by UniTester. Furthermore, the "Original+ChatUniTest" column reflects the total coverage resulting from the integration of the original tests with those generated by ChatUniTest. Finally, the "Original+RATester" column indicates the total coverage obtained by combining the original tests with those generated by RATester. Overall, both RATester and the baselines successfully enhance the line coverage of the original unit tests; however, RATester demonstrates a more significant improvement. Specifically, the line coverage increases substantially from a range of 0.51%–95.53% to 14.46%–95.57%. In addition, the average compile rate also shows a notable improvement, rising from 56.49% to 60.98%. This indicates that RATester not only enhances line coverage but also serves to complement human-written unit tests, thereby contributing to better overall software quality.

**Table 3.** RQ-1: The line coverage achieved by combining the original unit tests with those generated by RATester and the baselines

| Projects | Original | Original+ NxtUnit | Original+ UniTester | Original+ ChatUniTest | Original+ RATester |
|---|---|---|---|---|---|
| **beego** | 38.78% | 38.79% | 38.78% | 40.93% | 41.82% |
| **echo** | 93.58% | 93.58% | 93.59% | 93.68% | 93.96% |
| **fiber** | 85.46% | 86.15% | 85.56% | 86.35% | 86.59% |
| **frp** | 2.59% | 13.51% | 9.31% | 14.01% | 14.46% |
| **gin** | 95.53% | 95.54% | 95.54% | 95.56% | 95.57% |
| **hugo** | 76.54% | 76.58% | 76.60% | 76.77% | 77.20% |
| **nps** | 0.51% | 16.50% | 7.56% | 12.29% | 16.83% |
| **traefik** | 58.95% | 59.13% | 58.99% | 59.93% | 61.39% |
| **Average** | 56.49% | 59.97% | 58.24% | 59.94% | 60.98% |

**Effectiveness of RATester in Mutation Testing.** Table 4 presents the results of the study. For each project listed in column 1, the table details the number of generated mutants (column 2), the number of killed mutants by each approach (columns 3-6), and the mutator coverage percentages (columns 7-10). As shown in Table 4, we find that the unit tests generated by RATester not only kill the highest number of mutants but also achieve the best mutator coverage across all evaluated projects. For instance, in the "gin" project, a total of 885 mutants are generated. NxtUnit, UniTester, and ChatUniTest kill 127, 61, and 153 mutants, respectively. In contrast, RATester achieves an impressive 164 mutants killed, significantly surpassing the performance of the baselines. Furthermore, RATester achieves a mutator coverage of 26.05%, which is notably higher than NxtUnit's 18.21%, UniTester's 12.73%, and ChatUniTest's 25.76%. These results underscore the

effectiveness of RATester in not only detecting defects but also improving the overall quality and reliability of the generated unit tests when compared to existing approaches.

Table 4. RQ-1: RATester vs. Baselines across different projects in mutation testing

| Projects | Mutants (#) | Killed Mutants (#) | | | | Mutator Coverage (%) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | NxtUnit | UniTester | ChatUniTest | RATester | NxtUnit | UniTester | ChatUniTest | RATester |
| beego | 4,573 | 468 | 138 | 425 | 473 | 17.93% | 8.19% | 18.43% | 19.91% |
| echo | 922 | 49 | 26 | 90 | 91 | 15.59% | 12.21% | 16.39% | 16.59% |
| fiber | 1,554 | 153 | 80 | 141 | 160 | 15.24% | 9.36% | 15.93% | 17.99% |
| frp | 1,538 | 84 | 57 | 83 | 99 | 6.95% | 2.68% | 6.88% | 7.76% |
| gin | 885 | 127 | 61 | 153 | 164 | 18.21% | 12.73% | 25.76% | 26.05% |
| hugo | 5,882 | 542 | 339 | 620 | 670 | 9.27% | 7.75% | 11.07% | 11.36% |
| nps | 1,369 | 52 | 41 | 50 | 56 | 7.68% | 4.29% | 7.56% | 8.36% |
| traefik | 4,331 | 269 | 109 | 264 | 308 | 7.35% | 5.88% | 7.67% | 8.84% |
| Average | 2,632 | 218 | 106 | 228 | 253 | 12.28% | 7.89% | 13.71% | 14.61% |

> **Answer to RQ-1**: *RATester significantly outperforms baselines in enhancing compile rate and line coverage, improving from 16.67%−63.56% to 45.58%−69.49% and from 7.49%−53.92% to 12.92%−58.09%, respectively. It also surpasses other approaches in mutation testing, demonstrating its effectiveness in boosting software testing and quality.*

## 5.2 RQ-2: Model-Agnostic Capabilities of RATester

**Objective.** In RQ-1, we use DeepSeek-Coder as the backbone model to evaluate the effectiveness of RATester compared to the baselines. The results demonstrate that RATester outperforms existing approaches and shows promising performance in generating unit test cases. In this RQ, we extend our analysis to examine the model-agnostic effectiveness of RATester, specifically assessing whether RATester maintains its effectiveness when applied to different models.

**Experimental Design.** In addition to DeepSeek-Coder, we utilize two state-of-the-art open-source LLMs to investigate whether RATester remains effective across different models, thus evaluating its model-agnostic capabilities. Specifically, the additional models are (1) CodeLlama [35] and (2) Magicoder [46]. We follow the same experimental setup outlined in Section 4 and compare the performance of each model in terms of compile rate, line coverage, and mutation testing. To ensure consistency, we maintain identical experimental conditions across all basic LLMs and their corresponding RATester implementations.

**Results.** We discuss the model-agnostic capabilities of RATester from the aspects of compile rate, line coverage, and mutation testing, respectively.

**Model-agnostic capabilities of RATester in compile rate.** Fig. 4 shows the compile rate of unit tests generated by RATester compared to basic LLMs. We find that RATester significantly improves the performance of these basic LLMs. For instance, in the fiber project, RATester increases the compile rate from 44.8% to 61.8% for CodeLlama, from 37.4% to 59.6% for DeepSeek-Coder, and from 38.7% to 63.4% for Magicoder. Overall, RATester raises the compile rate of basic LLMs from a range of 30.6%-66.7% to 45.6%-74.4%, making more unit tests usable by developers during testing. This not only demonstrates the effectiveness of the RATester approach but also highlights its universal applicability. It is designed to be model-agnostic, meaning it can adapt to various LLMs, further emphasizing its flexibility and universality.

To further understand why RATester has an outstanding performance in unit test generation, we conduct a case study by analyzing one example (i.e., the unit tests generated by CodeLlama and RATester for the gin project) shown in Fig. 5. On the left of Fig. 5 is the focal method, the middle
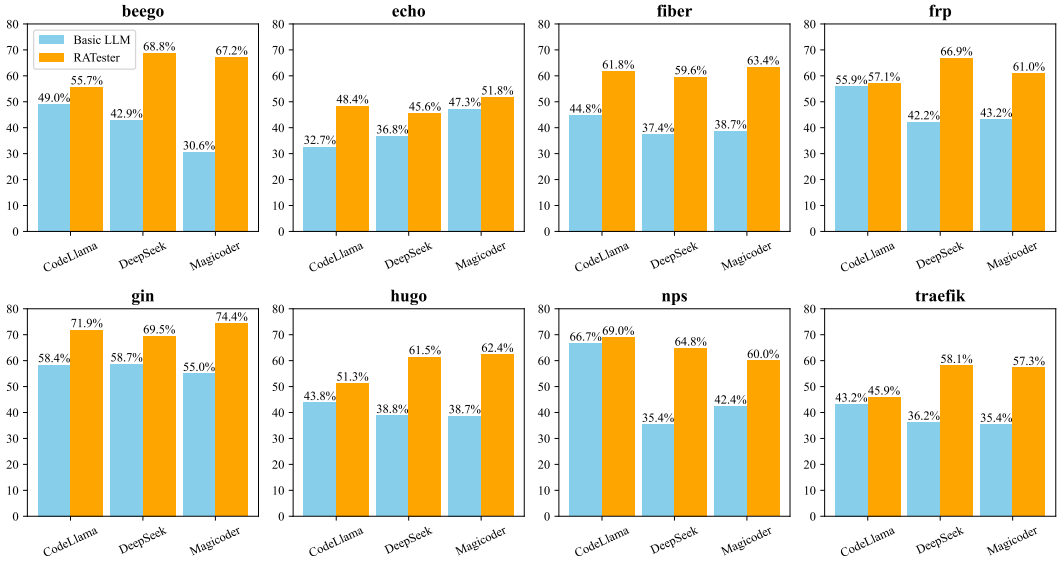
**Fig. 4.** RQ-2: RATester vs. Basic LLMs across different projects on compile rate

presents the unit test generated by CodeLlama, and the right shows the one generated by RATester. For the test case generated by CodeLlama, it tests the "Render" method of a struct named "String". The test creates an instance of the "String" type, sets the content type to "text/html" and the data to "Hello, World!". It then creates a mock HTTP response recorder using "httptest.NewRecorder", invokes the "Render" method to render the data, and checks if any errors occurred, if the HTTP status code is 200 OK, and if the response body correctly contains "Hello, World!". If an issue is detected at any checkpoint, the test reports an error using "t.Errorf" and terminates. However, this test case contains a compilation error, with the following runtime error: "cannot use 'Hello, World!' (untyped string constant) as []any value in struct literal". This occurs because CodeLlama hallucinates when generating the test case for the focal method, lacking knowledge about the "String" struct (i.e., no code context). As a result, it incorrectly infers the "Data" field and mistakenly sets it as a string. In contrast, RATester dynamically retrieves the required definitions and documentation when generating test cases. As a result, it correctly identifies that the Data field in the "String" struct should be a slice of any type, not a single string. This allows RATester to generate compilable unit test.

| Focal Method | Unit Test Generated by CodeLlama | Unit Test Generated by RATester |
|---|---|---|
| 01 ### Context_0 | func TestRender(t *testing.T) { | func TestRender(t *testing.T) { |
| 02 type String struct { | r := String{ | r := String{ |
| 03    Format string | Format: "text/html", | Format: "Hello, %s!", |
| 04    Data []any | Data: "Hello, World!", | Data: []any{"world"} |
| 05 } | } | } |
| 06 | w := httptest.NewRecorder() | w := httptest.NewRecorder() |
| 07 ### Context_1 | err := r.Render(w) | err := r.Render(w) |
| 08 ...... | if err != nil { | if err != nil { |
| 09 | t.Errorf("Render() error = %v", err) | t.Errorf("Render() error = %v", err) |
| 10 ### Focal Method | return | return |
| 11 func (r String) Render(w http.ResponseWriter) error | } | } |
| 12 { | if w.Code != http.StatusOK { | if w.Code != http.StatusOK { |
| 13    return WriteString(w, r.Format, r.Data) | ...... | ...... |
| 14 } | } | } |

**Fig. 5.** RQ-2: The unit tests generated by CodeLlama and RATester for the gin project

**Model-agnostic capabilities of RATester in line coverage.** The left side of Table 5 shows the line coverage results of basic LLMs, while the right side presents the results of RATester using different backbone models. We calculate not only the line coverage for RATester with different backbone models but also the improvement percentages relative to basic LLMs. As shown in the table, the RATester approach increases the overall line coverage of unit tests generated by basic LLMs across different projects. For example, in the beego project, RATester improves CodeLlama's coverage by 24.3% (from 20.95% to 26.04%), DeepSeek-Coder's coverage by 40.3% (from 22.75% to 31.91%), and Magicoder's coverage by 55.3% (from 22.35% to 34.71%). Overall, RATester enhances the performance of basic LLMs, raising their average line coverage from a range of 18.80%-19.95% to 23.00%-26.25%, with relative improvements ranging from 22.3% to 34.1%. This aligns with the motivation behind our method's design: providing LLMs with more effective context (such as definitions of called methods) enhances the model's ability to generate unit tests and subsequently increases overall line coverage.

**Table 5.** RQ-2: RATester vs. Basic LLMs across different projects on line coverage

| Projects | Basic LLM | | | RATester | | |
|---|---|---|---|---|---|---|
| | CodeLlama | DeepSeek | Magicoder | CodeLlama | DeepSeek | Magicoder |
| **beego** | 20.95% | 22.75% | 22.35% | 26.04% (↑24.3%) | 31.91% (↑40.3%) | 34.71% (↑55.3%) |
| **echo** | 17.65% | 21.08% | 24.21% | 26.35% (↑49.3%) | 24.50% (↑16.2%) | 27.55% (↑13.8%) |
| **fiber** | 14.94% | 17.64% | 14.82% | 20.72% (↑38.7%) | 26.31% (↑49.1%) | 16.41% (↑10.7%) |
| **frp** | 11.68% | 11.00% | 12.95% | 14.11% (↑20.8%) | 14.43% (↑31.2%) | 18.54% (↑43.2%) |
| **gin** | 43.53% | 45.35% | 42.67% | 48.23% (↑10.8%) | 58.09% (↑28.1%) | 47.13% (↑10.5%) |
| **hugo** | 16.10% | 16.87% | 16.02% | 19.77% (↑22.8%) | 25.02% (↑48.3%) | 18.92% (↑18.1%) |
| **nps** | 11.83% | 10.22% | 11.33% | 13.12% (↑10.9%) | 16.82% (↑64.6%) | 14.93% (↑31.8%) |
| **traefik** | 13.72% | 11.68% | 15.26% | 15.64% (↑14.0%) | 12.92% (↑10.6%) | 18.62% (↑22.0%) |
| **Average** | 18.80% | 19.57% | 19.95% | 23.00% (↑22.3%) | 26.25% (↑34.1%) | 24.60% (↑23.3%) |

**Model-agnostic capabilities of RATester in mutation testing.** Table 6 presents the results of RATester in using different backbone models in terms of mutation testing. From the results, we find that: **(1) Performance Variation Across Backbone Models:** There are significant performance differences among the basic LLMs, which directly impacts the effectiveness of RATester. Among the various backbone models tested, DeepSeek-Coder demonstrates superior performance, leading to the highest effectiveness when RATester utilizes DeepSeek-Coder as its backbone model. **(2) Enhancement Across All Models:** RATester consistently improves the performance of all three basic LLMs utilized in this study. For instance, in the hugo project, Magicoder kills only 522 mutants. In contrast, RATester using Magicoder successfully kills 583 mutants, showcasing a clear enhancement in defect detection capabilities. **(3) Overall Improvement in Defect Detection:** Across all three models tested, RATester exhibits a significant advantage, killing between 316 and 522 more mutants compared to the basic LLMs. This indicates that RATester not only leverages the strengths of the underlying models but also enhances their overall effectiveness in detecting defects.

**Answer to RQ-2**: *The basic LLMs have limited capabilities in generating unit tests, while RATester enhances these capabilities through appropriate adaptations. Overall, RATester significantly outperforms the basic LLMs in compile rate, line coverage, and mutation testing, demonstrating its adaptability and improved effectiveness.*

Table 6.  RQ-2: RATester vs. Basic LLMs across different projects in mutation testing

| Projects | Mutants (#) | Basic LLM | | | RATester | | |
|---|---|---|---|---|---|---|---|
| | | CodeLlama | DeepSeek | Magicoder | CodeLlama | DeepSeek | Magicoder |
| beego | 4,573 | 209 | 381 | 224 | 388 (+179) | 473 (+92) | 524 (+300) |
| echo | 922 | 42 | 48 | 50 | 48 (+6) | 91 (+43) | 61 (+11) |
| fiber | 1,554 | 126 | 131 | 115 | 131 (+5) | 160 (+29) | 124 (+9) |
| frp | 1,538 | 67 | 70 | 69 | 80 (+13) | 99 (+29) | 113 (+44) |
| gin | 885 | 59 | 135 | 86 | 78 (+19) | 164 (+29) | 115 (+29) |
| hugo | 5,882 | 428 | 615 | 522 | 448 (+20) | 670 (+55) | 583 (+61) |
| nps | 1,369 | 34 | 50 | 46 | 46 (+12) | 56 (+6) | 56 (+10) |
| traefik | 4,331 | 145 | 241 | 156 | 207 (+62) | 308 (+67) | 214 (+58) |
| Sum | 21,054 | 1,110 | 1,671 | 1,268 | 1,426 (+316) | 2,021 (+350) | 1,790 (+522) |

## 5.3  RQ-3: Efficiency of RATester

**Objective.** The previous research question examined the effectiveness of RATester. In this RQ, we aim to study the efficiency of RATester. We conduct a comprehensive experiment to evaluate its efficiency, focusing not only on the time required to generate test cases but also on the impact of the candidate number of generated unit tests for each focal method or function.

**Experimental Design.** We begin by investigating the total time required to generate test cases for all eight projects. We use the baselines (i.e., NxtUnit, UniTester, ChatUniTest, CodeLlama, DeepSeek-Coder, and Magicoder) to compare the total generation time of RATester across all projects. Additionally, we conduct a comparative analysis of the number of unit test candidates. We use DeepSeek-Coder and RATester (with DeepSeek-Coder as the backbone model) to generate test cases for all projects, setting the number of unit test candidates from 1 to 10 (i.e., generating 1 to 10 test cases for each focal method or function). We then calculate the average line coverage across all projects.

**Results.** We discuss the results from two perspectives: the total generation time across all projects and the impact of the unit test candidate number.

**Total generation time across all projects.** The left side of Fig. 6 shows the total generation time required by each baseline and RATester for unit test generation. According to the results, we observe that: (1) RATester requires more generation time compared to the basic LLMs. For example, CodeLlama alone takes 18.2 hours, while RATester built on CodeLlama requires 25.6 hours. Overall, the basic LLMs need between 18.2 and 19.9 hours, whereas RATester requires between 25.6 and 28.7 hours. However, we believe this additional time is acceptable in practical usage due to the higher compile rates, increased line coverage, and a greater number of killed mutants achieved by RATester (refer to Section 5.2 for more details). (2) The three basic LLMs (i.e., CodeLlama, DeepSeek-Coder, and Magicoder) perform fast, while the others perform relatively a little slow. More precisely, CodeLlama, DeepSeek-Coder, and Magicoder only take 18.2 hours, 19.1 hours, and 19.9 hours to generate unit tests for all projects. Among all the approaches, ChatUniTest has the longest generation time, taking 45.9 hours. This is because ChatUniTest includes redundant context and also features a unit test repair process that requires multiple iterations to arrive at the final unit test.

**Impact of the unit test candidate number.** According to the results on the right side of Fig. 6, we find that: (1) Different candidate numbers have varying impacts on the performance of RATester and DeepSeek-Coder, with both models showing improved performance as the number of candidates increases. (2) The curve indicates that the performance of RATester far exceeds that of DeepSeek-Coder. RATester only requires generating one unit test for each focal method or function to achieve

higher line coverage than DeepSeek-Coder, which requires ten unit tests. This demonstrates that RATester is more efficient and produces higher-quality unit tests with a less number of candidates. (3) Increasing the number of candidates does not guarantee significant performance improvements. As we continuously increase the number of candidates from 2 to 10, the performance of both RATester and DeepSeek-Coder improves only slightly, while the generation cost with the LLM increases significantly. Considering both the performance improvements and the associated generation costs, we adopt a candidate number of 1 unit test as the default setting.



**Fig. 6.** RQ-3: (Left) Total generation time of RATester and baselines; (Right) Performance of RATester and DeepSeek-Coder across varying unit test candidate number

> **Answer to RQ-3**: *(1) RATester requires more generation time than the basic LLMs; however, considering the performance improvements, this additional time is justified. (2) Increasing the candidate number can enhance the performance of RATester, but the improvement is not significant.*

## 6 Threats to Validity

**Internal Validity.** The first one arises from potential data leakage since referenced unit tests may be part of the training data of LLMs (e.g., CodeLlama and DeepSeek-Coder). To tackle this issue, we initially calculate the number of unit tests generated by RATester, which matches the reference unit test in all eight projects. We find that out of 11,195 generated unit tests, only 1 of these aligns with the unit tests in projects. Additionally, compared to the basic LLMs (i.e., CodeLlama, DeepSeek-Coder, and Magicoder), RATester demonstrates a significant enhancement in performance, achieving an increase in line coverage of 22.3% to 34.1%. Furthermore, the unit tests generated by RATester also contribute to completing the original unit tests in the projects, raising the line coverage from 56.49% to 60.98%. This demonstrates that the improved results achieved by RATester are not merely a result of memorizing the training data.

The second concern arises from potential errors in implementing our approach and baselines. To mitigate this threat, we implement our model through pair programming, using the source code of baselines provided by the corresponding authors and adhering to the same settings outlined in the original papers. Additionally, the authors conduct a thorough review of the experimental scripts to ensure their correctness.

**External Validity.** The main external threat to validity comes from our evaluation dataset used. The effectiveness demonstrated by RATester may not be generalizable to different unit test generation datasets, particularly those involving unit tests written in other programming languages (e.g., Java and Python). This limitation is common to some pipelines that utilize LLMs and language-specific tools (i.e., gopls), and we aim to address it in our future work.

## 7 Related Work

### 7.1 Unit Test Generation

Unit test generation approaches can be classified into three types: traditional approaches, learning-based approaches, and LLM-based approaches.

Traditional approaches [13, 29] focus on code coverage, and research shows that traditional approaches are very effective at achieving high coverage [7, 26, 30, 31]. Pacheco et al. introduced Randoop [29], a widely adopted tool that generates unit tests for Java. Fraser et al. developed EvoSuite [13], which automatically creates test suites optimized for high coverage, minimal size, and rich assertions. However, previous studies show that these approaches do not produce well-written, maintainable unit tests explicitly for developers to use [42, 44].

To address limitations in traditional unit test generation, learning-based approaches [6, 11, 18, 37, 42] have made notable advancements. Saes [37] leveraged over 780K focal-test method pairs from GitHub's JUnit framework to generate Java test suites with an 86.69% parsability rate. Tufano et al. [42] introduced AthenaTest, which fine-tunes the BART model on the Methods2Test dataset [41] to generate entire unit tests from focal method contexts, achieving a correct test rate for 43% of focal methods, with 16% classified as valid tests. Dinella et al. [11] proposed TOGA, which is a unified transformer-based neural approach to infer both exceptional and assertion test oracles based on the context of the focal method. Alagarsamy et al. [6] developed A3Test, which combines test oracle generation with domain adaptation to enhance naming consistency and signature verification, outperforming AthenaTest in test accuracy and method coverage. He et al. [18] proposed UniTester, which is trained on the UniTSyn dataset and is capable of synthesizing unit tests for programs in multiple languages, including Golang. However, these approaches rely heavily on task-specific datasets extracted from open-source repositories.

In response to the challenges posed by learning-based approaches, researchers are increasingly utilizing pre-trained LLMs to generate unit tests directly from contextual information, reducing reliance on task-specific datasets by leveraging extensive pre-training on diverse open-source code. Lemieux et al. [21] introduced CodaMOSA, an SBST approach that leverages LLMs to overcome coverage plateaus in Python code. Following that, Ni et al. [24] also explored ChatGPT for generating unit tests based on focal methods. Despite these advancements, LLMs may still exhibit hallucinations when generating unit tests for focal methods and functions due to their lack of the project's global knowledge. These hallucinations can include but are not limited to, calling non-existent methods, as well as assigning incorrect parameters and return values (e.g., mismatched parameter types or incorrect parameter counts). To overcome this limitation, many studies have explored the extraction of context to reduce hallucinations in the generation process of LLMs. Yuan et al. [54] developed ChatTester, an LLM-based model employing ChatGPT with an iterative generate-and-validate strategy that incorporates execution feedback and code context. Chen et al. [9] proposed ChatUniTest, which is an innovative framework designed to enhance automated unit test generation. It utilizes an LLM-based approach, augmented with an adaptive focal context mechanism to capture relevant context in prompts, and employs a "Generation-Validation-Repair" process to correct errors in generated tests. Following that, researchers [15, 36, 53] have explored the roles of focal context and dependency context. These methods utilize one or more fixed patterns to extract context for the focal method: (1) focal class signature; (2) signatures of other methods and fields in the class; (3) signatures of dependent classes; and (4) signatures of dependent methods and fields in the dependent classes. However, these fixed extraction patterns present several issues: (1) they may overlook important context; for instance, when generating a unit test for a specific focal method, the LLM might require unknown context beyond the dependencies of that focal

method; (2) there is potential for redundant context, as excessive irrelevant context could lead to redundancy, preventing the model from focusing on essential information.

Different from existing works, our paper presents an LLM-based framework that not only eliminates the reliance on task-specific datasets but also proactively utilizes gopls to fetch precise context. This approach significantly reduces hallucinations during the test generation process.

### 7.2 Pre-trained Language Model

With advancements in Natural Language Processing, Pre-trained Language Models have gained widespread traction due to their ability to be trained on billions of parameters and vast datasets, which has led to remarkable performance improvements across diverse applications. These models are highly adaptable to a range of downstream tasks, utilizing methods such as fine-tuning [33, 50] and prompting [23, 24, 49, 51]. Their versatility arises from extensive pre-training on broad data, equipping them with a robust knowledge base applicable across numerous domains. Fine-tuning involves adjusting model parameters specifically for a targeted task, requiring iterative training on a dedicated dataset, which enhances the model's accuracy and relevance for that task. By contrast, prompting offers a more direct, efficient approach by feeding the model task-specific instructions or a few relevant examples in natural language, allowing it to perform effectively without parameter adjustments. Although fine-tuning can yield higher accuracy, it demands significant computational resources and is less feasible in scenarios with limited task-specific datasets.

Pre-trained Language Models are typically based on the transformer architecture [43] and are categorized into three types: encoder-only, encoder-decoder, and decoder-only architectures. Encoder-only models, such as CodeBERT [12] and GraphCodeBERT [17], and encoder-decoder models, like PLBART [4] and CodeT5 [45], are trained using objectives like Masked Language Modeling (MLM) or Masked Span Prediction (MSP). In these setups, a small percentage (e.g., 15%) of the tokens are replaced with either masked tokens or masked span tokens, and the model learns to predict or recover the masked content. Trained on diverse code-related data, these models are then fine-tuned for specific tasks to achieve enhanced performance [14, 19, 25]. Decoder-only models have gained significant attention, primarily due to their use of causal language modeling objectives, which train them to predict the probability of each next token based on all previous tokens in a sequence. GPT [33] and its variants are the most prominent examples of this architecture, marking a pivotal point in bringing large language models into widespread practical applications.

To enhance LLMs' generalization and alignment with human intentions on previously unseen downstream tasks, recent research has focused on instruction tuning and reinforcement learning to improve model performance [10, 28, 55]. For instance, OpenAI's ChatGPT [27] is a notable example built on the generative pre-trained transformer architecture. It undergoes initial instruction tuning, followed by updates through reinforcement learning from human feedback to better capture human-aligned responses. Beyond commercial models, there is also a growing landscape of open-source instructed LLMs, such as CodeLlama [35] and DeepSeek-Coder [5], which demonstrate promising performance across various tasks and hold potential for broader adaptability [24, 50, 52].

### 8 Conclusion

This paper enhances the LLM's ability to generate more repository-aware unit tests through global contextual information injection. To provide LLMs with a level of global knowledge similar to that of human testers, RATester integrates the language server gopls. When it encounters unfamiliar identifiers, such as struct names, RATester utilizes gopls to fetch relevant precise context, thereby preventing erroneous usage. This integration enriches the LLM's global knowledge of the project, significantly reducing hallucinations. We evaluate the effectiveness and efficiency of RATester by constructing a new Golang dataset from real-world projects and comparing it against baseline

approaches. The results illustrate the advantages of RATester over these baselines. Furthermore, we extend our analysis to assess the model-agnostic effectiveness of RATester. These findings not only validate the efficacy of RATester but also emphasize its universal applicability.

## Acknowledgements

## References

[1]  2024. gopls.  https://github.com/golang/tools/tree/master/gopls
[2]  2024. gremlins.  https://github.com/go-gremlins/gremlins
[3]  2024. Hugging Face.  https://huggingface.co
[4]  Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Unified pre-training for program understanding and generation. *arXiv preprint arXiv:2103.06333* (2021).
[5]  DeepSeek AI. 2023. DeepSeek Coder: Let the Code Write Itself. https://github.com/deepseek-ai/DeepSeek-Coder.
[6]  Saranya Alagarsamy, Chakkrit Tantithamthavorn, and Aldeida Aleti. 2023. A3Test: Assertion-Augmented Automated Test Case Generation. *arXiv preprint arXiv:2302.10352* (2023).
[7]  Aldeida Aleti, Irene Moser, and Lars Grunske. 2017. Analysing the fitness landscape of search-based software testing problems. *Automated Software Engineering* 24 (2017), 603–621.
[8]  Earl T Barr, Mark Harman, Phil McMinn, Muzammil Shahbaz, and Shin Yoo. 2014. The oracle problem in software testing: A survey. *IEEE transactions on software engineering* 41, 5 (2014), 507–525.
[9]  Yinghao Chen, Zehao Hu, Chen Zhi, Junxiao Han, Shuiguang Deng, and Jianwei Yin. 2024. Chatunitest: A framework for llm-based test generation. In *Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering*. 572–576.
[10]  Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems* 30 (2017).
[11]  Elizabeth Dinella, Gabriel Ryan, Todd Mytkowicz, and Shuvendu K Lahiri. 2022. Toga: A neural method for test oracle generation. In *Proceedings of the 44th International Conference on Software Engineering*. 2130–2141.
[12]  Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. Codebert: A pre-trained model for programming and natural languages. *arXiv preprint arXiv:2002.08155* (2020), 1536–1547.
[13]  Gordon Fraser and Andrea Arcuri. 2011. 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*. 416–419.
[14]  Michael Fu and Chakkrit Tantithamthavorn. 2022. Linevul: A transformer-based line-level vulnerability prediction. In *Proceedings of the 19th International Conference on Mining Software Repositories*. 608–620.
[15]  Shuzheng Gao, Chaozheng Wang, Cuiyun Gao, Xiaoqian Jiao, Chun Yong Chong, Shan Gao, and Michael Lyu. 2025. The Prompt Alchemist: Automated LLM-Tailored Prompt Optimization for Test Case Generation. *arXiv preprint arXiv:2501.01329* (2025).
[16]  Vahid Garousi and Junji Zhi. 2013. A survey of software testing practices in Canada. *Journal of Systems and Software* 86, 5 (2013), 1354–1376.
[17]  Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, et al. 2020. Graphcodebert: Pre-training code representations with data flow. *arXiv preprint arXiv:2009.08366* (2020).
[18]  Yifeng He, Jiabo Huang, Yuyang Rong, Yiwen Guo, Ethan Wang, and Hao Chen. 2024. UniTSyn: A Large-Scale Dataset Capable of Enhancing the Prowess of Large Language Models for Program Testing. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*. 1061–1072.
[19]  David Hin, Andrey Kan, Huaming Chen, and M Ali Babar. 2022. LineVD: Statement-level Vulnerability Detection using Graph Neural Networks. *arXiv preprint arXiv:2203.05181* (2022).
[20]  Jihyun Lee, Sungwon Kang, and Danhyung Lee. 2012. Survey on software testing practices. *IET software* 6, 3 (2012), 275–282.
[21]  Caroline Lemieux, Jeevana Priya Inala, Shuvendu K Lahiri, and Siddhartha Sen. 2023. 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, 919–931.

[22] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023. StarCoder: may the source be with you! *arXiv preprint arXiv:2305.06161* (2023).

[23] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *Comput. Surveys* 55, 9 (2023), 1–35.

[24] Chao Ni, Xiaoya Wang, Liushan Chen, Dehai Zhao, Zhengong Cai, Shaohua Wang, and Xiaohu Yang. 2024. CasModaTest: A Cascaded and Model-agnostic Self-directed Framework for Unit Test Generation. *arXiv preprint arXiv:2406.15743* (2024).

[25] Chao Ni, Xin Yin, Kaiwen Yang, Dehai Zhao, Zhenchang Xing, and Xin Xia. 2023. Distinguishing Look-Alike Innocent and Vulnerable Code by Subtle Semantic Representation Learning and Explanation. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1611–1622.

[26] Carlos Oliveira, Aldeida Aleti, Lars Grunske, and Kate Smith-Miles. 2018. Mapping the effectiveness of automated test suite generation techniques. *IEEE Transactions on Reliability* 67, 3 (2018), 771–785.

[27] OpenAI. 2022. ChatGPT: Optimizing Language Models for Dialogue. (2022). https://openai.com/blog/chatgpt/.

[28] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* 35 (2022), 27730–27744.

[29] Carlos Pacheco and Michael D Ernst. 2007. Randoop: feedback-directed random testing for Java. In *Companion to the 22nd ACM SIGPLAN conference on Object-oriented programming systems and applications companion*. 815–816.

[30] Annibale Panichella, Fitsum Meshesha Kifetew, and Paolo Tonella. 2015. Reformulating branch coverage as a many-objective optimization problem. In *2015 IEEE 8th international conference on software testing, verification and validation (ICST)*. IEEE, 1–10.

[31] Annibale Panichella, Fitsum Meshesha Kifetew, and Paolo Tonella. 2017. Automated test case generation as a many-objective optimisation problem with dynamic selection of the targets. *IEEE Transactions on Software Engineering* 44, 2 (2017), 122–158.

[32] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 8024–8035. http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf

[33] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).

[34] Nikitha Rao, Kush Jain, Uri Alon, Claire Le Goues, and Vincent J Hellendoorn. 2023. CAT-LM training language models on aligned code and tests. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 409–420.

[35] Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950* (2023).

[36] Gabriel Ryan, Siddhartha Jain, Mingyue Shang, Shiqi Wang, Xiaofei Ma, Murali Krishna Ramanathan, and Baishakhi Ray. 2024. Code-aware prompting: A study of coverage-guided test generation in regression setting using llm. *Proceedings of the ACM on Software Engineering* 1, FSE (2024), 951–971.

[37] Laurence Saes. 2018. Unit test generation using machine learning. *Universiteit van Amsterdamg* (2018).

[38] Max Schäfer, Sarah Nadi, Aryaz Eghbali, and Frank Tip. 2023. An empirical evaluation of using large language models for automated unit test generation. *IEEE Transactions on Software Engineering* (2023).

[39] Jessica Shieh. 2023. Best practices for prompt engineering with OpenAI API. *OpenAI, February https://help.openai.com/en/articles/6654000-best-practices-for-prompt-engineering-with-openai-api* (2023).

[40] Jiho Shin, Sepehr Hashtroudi, Hadi Hemmati, and Song Wang. 2024. Domain Adaptation for Code Model-Based Unit Test Case Generation. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*. 1211–1222.

[41] Michele Tufano, Shao Kun Deng, Neel Sundaresan, and Alexey Svyatkovskiy. 2022. Methods2Test: A dataset of focal methods mapped to test cases. In *Proceedings of the 19th International Conference on Mining Software Repositories*. 299–303.

[42] Michele Tufano, Dawn Drain, Alexey Svyatkovskiy, Shao Kun Deng, and Neel Sundaresan. 2020. Unit test case generation with transformers and focal context. *arXiv preprint arXiv:2009.05617* (2020).

[43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).

[44] Siwei Wang, Xue Mao, Ziguang Cao, Yujun Gao, Qucheng Shen, and Chao Peng. 2023. NxtUnit: Automated Unit Test Generation for Go. In *Proceedings of the 27th International Conference on Evaluation and Assessment in Software Engineering*. 176–179.

[45] Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. 2021. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. *arXiv preprint arXiv:2109.00859* (2021).

[46] Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. 2023. Magicoder: Source code is all you need. *arXiv preprint arXiv:2312.02120* (2023).

[47] Chunqiu Steven Xia, Yuxiang Wei, and Lingming Zhang. 2023. Automated program repair in the era of large pre-trained language models. In *Proceedings of the 45th International Conference on Software Engineering (ICSE 2023). Association for Computing Machinery*.

[48] Chunqiu Steven Xia and Lingming Zhang. 2023. Keep the Conversation Going: Fixing 162 out of 337 bugs for $0.42 each using ChatGPT. *arXiv preprint arXiv:2304.00385* (2023).

[49] Xin Yin, Chao Ni, Tien N Nguyen, Shaohua Wang, and Xiaohu Yang. 2024. Rectifier: Code Translation with Corrector via LLMs. *arXiv preprint arXiv:2407.07472* (2024).

[50] Xin Yin, Chao Ni, and Shaohua Wang. 2024. Multitask-based evaluation of open-source llm on software vulnerability. *IEEE Transactions on Software Engineering* (2024).

[51] Xin Yin, Chao Ni, Shaohua Wang, Zhenhao Li, Limin Zeng, and Xiaohu Yang. 2024. Thinkrepair: Self-directed automated program repair. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*. 1274–1286.

[52] Xin Yin, Chao Ni, Xiaodan Xu, and Xiaohu Yang. 2024. What You See Is What You Get: Attention-based Self-guided Automatic Unit Test Generation. *arXiv preprint arXiv:2412.00828* (2024).

[53] Zhiqiang Yuan, Mingwei Liu, Shiji Ding, Kaixin Wang, Yixuan Chen, Xin Peng, and Yiling Lou. 2024. Evaluating and improving chatgpt for unit test generation. *Proceedings of the ACM on Software Engineering* 1, FSE (2024), 1703–1726.

[54] Zhiqiang Yuan, Yiling Lou, Mingwei Liu, Shiji Ding, Kaixin Wang, Yixuan Chen, and Xin Peng. 2023. No more manual tests? evaluating and improving chatgpt for unit test generation. *arXiv preprint arXiv:2305.04207* (2023).

[55] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593* (2019).