

LSPAI: An IDE Plugin for LLM-Powered Multi-Language Unit Test Generation with Language Server Protocol

Gwihwan Go

BNRist, Tsinghua University
Beijing, China
iejw1914@gmail.com

Chijin Zhou

BNRist, Tsinghua University
Beijing, China
tlock.chijin@gmail.com

Quan Zhang

BNRist, Tsinghua University
Beijing, China
zhangq20@mails.tsinghua.edu.cn

Yu Jiang*

BNRist, Tsinghua University
Beijing, China
jiangyu198964@126.com

Zhao Wei

Tencent
Beijing, China
zachwei@tencent.com

Abstract

Unit testing is crucial for ensuring code validity, and extensive research has been conducted to advance this domain. However, existing studies fail to address critical industry requirements, particularly support for multi-language static analysis and real-time unit test generation. While integrating static analysis with a Large Language Model could address these challenges, it typically requires significant manual effort to implement across diverse programming languages. To address this, we propose LspAi, an automated unit test generation tool that leverages well-established language analysis tools and integrates them into a unified development environment via the Language Server Protocol. This approach equips LLM with multi-language static analysis capabilities, allowing a single tool to support systematic unit test generation across multiple languages. We evaluate our method by comparing line coverage across various LLMs and programming languages, demonstrating superior performance and broader applicability. Our evaluation of LspAi on real-world projects showed line coverage improvements of 114.3% for Java, 850% for Golang, and 3.03% for Python, along with enhancements in the valid ratio for most cases. In addition, we also share our lessons learned from applying the tool in Tencent Ltd.

ACM Reference Format:

Gwihwan Go, Chijin Zhou, Quan Zhang, Yu Jiang, and Zhao Wei. 2025. LSPAI: An IDE Plugin for LLM-Powered Multi-Language Unit Test Generation with Language Server Protocol. In *Companion Proceedings of the 33rd ACM Symposium on the Foundations of Software Engineering (FSE-Companion'25)*, June 23–27, 2025, Trondheim, Norway. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

Unit testing plays a pivotal role in modern software development by ensuring the validity and reliability of code. As software systems

grow in complexity, the importance of robust unit tests cannot be overstated, serving as a fundamental practice for identifying defects early and facilitating maintainable codebases. Extensive research has been dedicated to automating unit test generation, leading to the development of Search-Based Software Testing (SBST) tools such as EvoSuite [10], Randoop [25], and Pyguint [20]. More recently, the evolution of Large Language Models (LLMs) has introduced a new paradigm for unit test generation. Models like GPT [5] and Copilot [12] can understand the context of the code and generate semantically relevant unit tests, significantly improving software development efficiency.

Despite significant advancements, LLMs are still prone to generating incorrect unit tests. For example, Copilot, one of the most popular tools used by many companies, is still capable of making mistakes, as acknowledged by the Copilot development team [21]. Similarly, Siddiq et al. [28] found that LLM-generated test cases yield a relatively low validity rate, ranging from 2% to 12.7%, based on the SF110 [11] benchmark. As a result, researchers have proposed integrating static analysis with LLMs to help them better understand context and generate more accurate unit tests [16, 33, 35]. However, current research does not address the following two fundamental requirements of the software industry, which limits the broader adoption of these approaches in real-world settings.

First, performing static analysis across multiple programming languages is challenging. Industries adopt a variety of programming languages for different projects, and a test case generator should ideally support multiple languages. However, as shown in Table 1, most academic research has focused on one specific language rather than multi-language support. This focus stems from the difficulty of performing unified static analysis across diverse languages. Therefore, developers are forced to build customized analysis pipelines for each language, which requires significant manual adaptation effort. As a result, as far as we know, there are currently no academic tools available that can generate multi-language unit tests using static analysis.

Second, it is challenging to support real-time unit test generation when integrating static analysis. Developers often write unit tests concurrently with writing code. However, current SBST tools and LLM-integrated tools are unsuitable for scenarios requiring instant test generation, as they typically require the compilation of entire projects to perform static analysis and collect

*Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
Conference'17, July 2017, Washington, DC, USA

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

coverage feedback. Consequently, the reliance of SBST tools on coverage feedback to enhance test case quality limits their feasibility for real-time use. This issue also persists in recent LLM-integrated tools [1, 2, 16, 17, 33, 35], which depend on heavy static analysis and coverage feedback to mitigate LLM hallucinations. As a result, no academic tool currently supports real-time unit test generation, as illustrated in Table 1.

To tackle the aforementioned challenges, we introduce LSPAI, a real-time unit test generation tool powered by LLMs and integrated with static analysis for multi-language codebases. Our key insight is that well-established language analysis tools exist for each programming language and can be accessed through the Language Server Protocol (LSP) in a unified way. By leveraging the LSP, we can perform lightweight static analysis on multiple languages within a single environment with minimal effort. Specifically, LSPAI operates in two main steps: First, LSPAI conducts dependency analysis by extracting key tokens from the focal method and retrieves the corresponding dependent source code. Second, using the retrieved dependency source code, LSPAI performs real-time unit test generation and fixing. This approach effectively leverages reliable static analysis tools to improve LLMs' ability.

LSPAI brings two main benefits to developers who work in industries. First, LSPAI supports real-time unit test generation *without whole project compilation*, allowing developers to generate unit tests concurrently with code writing. Second, LSPAI simplifies the setup process by only requiring *a simple installation* of relevant language analysis plugins (e.g., extensions for Visual Studio Code), making it easily adaptable to various programming languages.

We developed LSPAI as an IDE (Integrated Development Environment) plugin for seamless integration and evaluated its performance across three widely used programming languages: Java, Python, and Golang. Our evaluation shows that LSPAI consistently improves unit test performance in terms of line coverage across real-world projects, regardless of the programming language. Specifically, it achieves a line coverage increase of 181.78% for Java, 850% for Golang, and 3.03% for Python compared to the baseline. Additionally, LSPAI improves the valid ratio for the target real-world projects, achieving increases of 114.39% for Java and 483.6% for Golang and 3.90% for Python. In addition, we also share our lessons learned from applying the tool in Tencent Ltd.

- We identified a research gap in current unit test generation: the lack of support for multi-language codebases and real-time test generation scenarios.
- We designed LSPAI as an IDE plugin, a practical tool that generates effective unit tests across multiple programming languages.

Table 1: Literature analysis which shows current research gap on unit test generation.

Tools	Real-Time	Multi-Language	Static Analysis Support				
			Java	Python	Golang	Others	
UTGen [4], EvoSuite [10, 18, 37], Randoop [25], HITS [33], casmoda [24], testspark [27], ChatUnitTest [35]	✗	✗	✓	✗	✗	✗	
PynGuin [20], CodaMosa [16], CoverUp [1], MuTAP [3], TELPA [36], SymPrompt [26], CLAP [31]	✗	✗	✗	✓	✗	✗	
NxtUnitGo [32]	✗	✗	✗	✗	✓	✗	
Copilot [12]	✓	✓	✗	✗	✗	✗	
LSPAI	✓	✓	✓	✓	✓	✓	

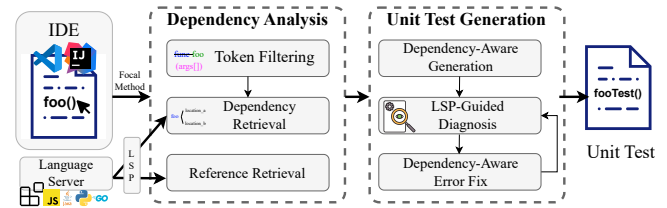


Figure 1: Overall workflow of LSPAI.

The source code can be found at <https://github.com/Gwihwan-Go/LSPAI/tree/fse-industry>.

- We evaluated LSPAI in real-world projects written in three programming languages, demonstrating its ability to consistently improve unit test performance. LSPAI achieved line coverage improvements of 114.3% for Java, 850% for Golang, and 3.03% for Python, along with enhanced valid test ratios for most projects.

2 Background

The *Language Server Protocol* (LSP) [22] is a standardized way for development tools, like text editors and IDEs, to communicate with *language servers*. A *Language Server* is meant to provide language-specific tasks, such as static analysis and code action recommendations. The main idea behind LSP is to provide a single, unified protocol that allows a language server to be used across different development environments, supporting multiple programming languages with minimal additional setup. This means that development tools can access advanced features for many programming languages through the same protocol, making it easier to work with different languages without needing to implement those features from scratch. Before LSP, features like syntax highlighting and code completion had to be written separately for each development environment and programming language. With LSP, however, editors can simply connect to a language server, which provides these features automatically, saving time and effort.

3 Design of LSPAI

This section describes the design of LSPAI, a unit test generation tool that enhances unit test creation through multi-language static analysis aided by LSP. Figure 1 illustrates LSPAI's overall workflow. When the developer requests unit test generation for a specific method, LSPAI generates a unit test following two steps. First, LSPAI collects dependency information for the given method by extracting and retrieving token definitions. Second, it generates the unit test based on the collected dependency information. The generated unit test is then analyzed using LSP. If any issues are detected, LSPAI retrieves the necessary dependencies and corrects the errors.

3.1 Employed LSP Features

LSPAI leverages the advantages of LSP, conducting static analysis using features provided by the language server. In this way, LSPAI can do consistent analysis across different IDEs and programming languages with the unified pipeline. Specifically, LSPAI utilizes the following features provided by the language server.

Symbol Provider. In LSP, symbols represent code entities such as files, modules, classes, functions, and variables within the source code. When LSPAI requests symbols for a specific file path, the

language server returns a hierarchical structure of these symbols found in the given text document. LSPA leverages the *Symbol Provider* to identify unit test entries through these symbols and to locate variable definitions by traversing the collected symbols.

Semantic Token Provider. Tokens are the smallest elements of code that can be broken down [23]. Semantic tokens extend tokens by adding contextual information, utilizing language servers that deeply understand the source file. When LSPA requests semantic tokens for a specific range, the language server returns an array of objects, each containing the context information of the token and its range. The *Semantic Token Provider* allows LSPA to determine how each token is used, facilitating a more granular and precise analysis of methods, variables, and other constructs.

Definition Provider. This feature plays a crucial role in analyzing dependencies for a target function by identifying the locations of function or class declarations. When LSPA requests the *Definition Provider* for a specific token, it returns the locations of token definitions. This capability allows LSPA to accurately map dependencies within the codebase, ensuring comprehensive analysis of the target functions or classes.

Reference Provider. The *Reference Provider* enables LSPA to identify all occurrences of a particular symbol within the codebase. By requesting references for a specific symbol, the language server returns a list of locations where the symbol is used. This feature is essential for understanding the context and usage patterns of the symbol, which aids in accurately determining dependencies and ensuring that generated unit tests cover relevant interactions.

Diagnosis Provider. It is vital for increasing the validity of the generated code. When LSPA requests a diagnosis, the *Diagnosis Provider* provided by the language server analyzes the code using its understanding of the source, detects any warnings or errors, and provides their specific locations. This allows LSPA to effectively identify and rectify issues in the generated code.

3.2 Dependency Analysis

This module gathers dependency information to generate reliable unit tests with high coverage. The dependency information is collected in three steps: token filtering, dependency retrieval, and reference retrieval. Through this process, LSPA acquires streamlined, high-quality dependency information.

Token Filtering. This step enhances the quality of the data to be collected by extracting tokens that are more likely to be relevant to the focal method. A focal method that requires testing is typically complex, containing numerous tokens. Analyzing every token of the method would generate a large amount of unnecessary data, most of which would not contribute to unit test generation. For example, the `parse` method in the `Parser` class of the `commons-cli` [6] project contains over 100 tokens within approximately 40 lines of code. Analyzing and retrieving information for over 100 tokens is inefficient and does not effectively enhance the quality of unit tests. Therefore, appropriate token filtering is essential. The token filtering strategy of LSPA involves two main steps: (1) *Selecting Key Tokens*: LSPA consider a token is important if it is given by the argument value of the method or is returned by the method. (2) *Selecting Associated Tokens*: LSPA determine whether the tokens are associated with key tokens, utilizing the

knowledge of the language server. Specifically, it requests the role of tokens that co-located with the key tokens by examining their types and modifiers through *Semantic Token Provider*. A co-located token is considered meaningful if the language server considers the token's role as declaring or defining. Following the above steps, the 100 tokens under `parse` method can be streamlined to 10 tokens. Ultimately, this module returns the extracted tokens, which LSPA uses to retrieve further information.

Dependency Retrieval. This step involves a strategy to extract relevant dependency information from the given tokens, ensuring LSPA retains only essential data for unit test generation while discarding unnecessary details from the language server. This is important because retrieved dependency information is often verbose, including comments, unrelated properties, or large code snippets. This can hinder unit test generation and degrade LSPA's performance. To address this, we apply heuristic rules based on LSP knowledge. First, by requesting the *Definition Provider* using the token's position, we collect the symbol that defines the token. Next, using the *Symbol Provider*, we identify the symbol's type (e.g., function, class, method, variable, or property). Finally, we summarize the relevant code snippet based on the symbol type. For example, functions are summarized by their return type and input arguments, while methods are summarized with their return type, input arguments, and associated class member properties.

Reference Retrieval. Referring to the use case of the focal method can enhance the correctness of generated test codes. Especially for LLM, which determines its output based on context, much research [13, 34] has proved that giving an example can enhance the quality of its output. In this regard, LSPA collects every use case of the focal method, utilizing *Reference Provider*. The collected reference information is then passed to the next step along with the dependency information and is used to generate unit test code.

3.3 Unit Test Generation

This module is responsible for generating reliable unit tests with high coverage without compiling or executing code. It leverages the given dependency information and LLM to generate unit tests. Subsequently, to mitigate the limitations of LLM, it detects issues in the generated code and fixes them.

Dependency-Aware Generation. LSPA generates unit tests by incorporating information passed by the section 3.2. In detail, We construct the prompt incorporating the source code of the focal method, natural language description, and retrieved information.¹ We construct our prompt template based on that of `ChatUniTest` [35]. Since this template is Java-specific, for generating unit tests in other programming languages, we slightly modify the prompt accordingly. The final prompt ranges from 1000 to 1,500 tokens, depending on the length of the focal method. The constructed prompt is then provided to a LLM to produce the unit test.

LSP-Guided Diagnosis. This component detects issues in the generated test code in real-time without the need for compilation or execution. LLMs can produce syntactically incorrect or incompatible code due to hallucinations. Hallucination [14, 15] refers to the generation of syntactically incorrect or semantically invalid

¹For detailed prompt snippets used by LSPA, refer to <https://github.com/GwihwanGo/LSPA/tree/fse-industry/src/promptBuilder.ts>

code that deviates from the desired output. However, in a real-time generation setting, where compilation or execution is not feasible, we need alternative methods to mitigate the LLM’s hallucinations. LSPAI utilizes *Diagnosis Provider* supported by LSP to inspect the generated code. If *Diagnosis Provider* does not detect any issues, LSPAI saves the unit test code, if there is any issue, LSPAI collects them and prioritize based on severity.

Dependency-Aware Error Fix. This step integrates dependency information to effectively fix errors. Based on the diagnosis, LSPAI identifies the related symbol and retrieves the necessary dependency information using the *Symbol Provider*. This information is then incorporated into the prompt to assist the LLM in correcting the error alongside the necessary dependency source code. The constructed prompt is sent to the LLM to fix the code. After the fix is made, LSPAI returns to the *LSP-Guided Diagnosis* step to verify whether the issue has been resolved. If the error is fixed or the iteration limit is reached, the corrected code is saved and presented to the developers.

4 Evaluation

In this section, we comprehensively evaluate the performance of LSPAI by evaluating the LSPAI’s performance on real-world projects across different programming languages.

4.1 Experiment Setup

Programming languages. We selected three different programming languages, Python, Java, and Golang, for the experiment. We selected Python and Java because their unit test generation capabilities have been extensively studied in previous research. Golang was chosen for two reasons: (1) it is widely used in industry but has received relatively little attention in academic studies, and (2) as a relatively new language, we anticipate it presents unique challenges for LLMs in generating valid code.

Real-world Projects. For a fair evaluation, we selected real-world projects from either (1) benchmarks used in previous research or (2) projects that are widely adopted by the community. As Table 2 shows, we selected two projects for each programming language. For Java code bases, we choose Commons-CLI [6] and Commons-CSV [7], which are commonly selected as benchmarks by previous research [33, 35] of unit test domain. For Golang code bases, we chose Logrus [29], and Cobra [9], because they are also selected by previous research [32] for evaluating unit test performance. Finally, for Python, we chose Black [8], which is also widely selected as a benchmark by previous research [16, 20]. We chose Crawl4AI [30] to test the ability of LLM’s code generation ability against a relatively new project that is not included in LLM training datasets. The Crawl4AI project is one of the most trending projects in GitHub with 24.6k stars, released in May 2024, which is well-developed while surely unseen in the training dataset of LLMs.

LLMs and Environment. We evaluate the effectiveness of LSPAI on three kinds of LLMs, including GPT4o [5], GPT4o-mini, and DeepSeek-V3 (DS-V3) [19], which are popular choices in LLM applications. We adopt the default temperature and generation settings for all LLMs. We conducted our evaluation on a machine equipped with an AMD EPYC 7763 CPU (2.25GHz) with 128 cores and 8 NVIDIA GPU (V100-32G), running Ubuntu 22.04 LTS.

Table 2: Dataset Statistics

Project	Abbr.	Domain	Version	Language
Commons-CLI [6]	CLI	Cmd-line Interface	eb541428	Java
Commons-CSV [7]	CSV	Csv file Processing	92e486ac	Java
Logrus [29]	LOG	logging for Golang	d1e6332	Golang
Cobra [9]	COB	Golang CLI interactions	3a6873e	Golang
Black [8]	BAK	Python code formatter	8dc9127	Python
Crawl4AI [30]	C4AI	LLM Friendly Web Crawler	8878b3d	Python

Table 3: Comparative experiment results on line coverage and valid rate.

	Model	Line Coverage		Valid Rate	
		LSPAI	NAIVE	LSPAI	NAIVE
CLI	GPT4o	57.04%	35.14%	56.52%	40.58%
	GPT4o-mini	51.34%	13.30%	43.00%	17.87%
	DS-V3	54.31%	28.72%	57.48%	44.93%
CSV	GPT4o	50.53%	35.62%	55.71%	31.43%
	GPT4o-mini	42.25%	17.56%	38.57%	6.43%
	DS-V3	68.26%	41.01%	43.57%	10.71%
LOG	GPT4o	42.17%	1.53%	30.00%	2.85%
	GPT4o-mini	33.16%	4.76%	11.42%	3.57%
	DS-V3	29.59%	8.50%	21.42%	14.28%
COB	GPT4o	11.24%	1.23%	34.64%	5.23%
	GPT4o-mini	10.72%	5.09%	16.99%	6.53%
	DS-V3	14.05%	1.81%	41.17%	27.45%
BAK	GPT4o	47.49%	44.79%	55.00%	43.18%
	GPT4o-mini	35.43%	33.89%	52.04%	59.54%
	DS-V3	38.04%	37.29%	71.36%	67.27%
C4AI	GPT4o	32.74%	32.72%	54.37%	49.33%
	GPT4o-mini	37.36%	36.28%	52.78%	62.59%
	DS-V3	40.70%	39.69%	67.90%	62.86%
Total		38.69%	23.20%	44.66%	30.92%

4.2 Comparative Experiment

This section evaluates the performance of LSPAI using two metrics: line coverage and valid rate, where a script is considered valid if it can be executed without errors. Since there are no exist open-sourced multi-language unit test generation tool, we constructed the baseline with the same prompt template with LSPAI, except for the dependency information collected from LSP. We refer to this baseline as NAIVE. As shown in Table 3, LSPAI significantly improves both line coverage and valid rate across three programming languages, different projects, and various LLMs. We will look into a more detailed analysis to find out what contributes to these results. Since there are slightly different trends in each programming language, we will examine each language individually.

Java. On average, LSPAI improves its baseline NAIVE by 114.39% in line coverage and 181.78% in valid rate. The primary reason for the increased coverage is the dependency retrieval-guided unit test generation, which effectively utilizes summarized dependency information of classes and methods to cover diverse edge cases. For the valid rate, LSPAI achieves substantial improvements by 181.78%

due to Java’s highly structured nature, which allows most errors to be detected before compilation through LSP. This is particularly evident in the substantial valid rate improvement observed in Java projects (e.g., CLI and CSV) compared to NAIVE.

Golang. On average, LSPAI enhances its baseline NAIVE by 850% in line coverage and 483.6% in the valid rate. For Golang-specific tasks, we found that LLMs often generate invalid unit test code due to simple mistakes, which is reflected in the lowest valid rate of NAIVE for Golang projects (e.g., LOG and CCOB). LSPAI addresses this by detecting and correcting these errors, leading to a significant improvement in the valid rate, which in turn boosts line coverage. For instance, during experiments, we observed that LLMs frequently “redeclare” objects that are already defined in the source code, causing invalid test code and resulting in low valid rates for the LOG and CCOB projects (averaging 6.89% and 13.07%). By leveraging LSP, LSPAI can detect and resolve this issue, effectively mitigating the limitations of LLMs. The rapid improvement in both line coverage and valid rate demonstrates that LSPAI can significantly enhance the quality of generated test code, particularly for programming languages where LLMs tend to struggle.

Python. For Python projects such as BAK and C4AI, LSPAI achieved an averaged modest increase in line coverage of 3.03%, and a increase in the valid rate of 3.90% on average. These relatively low improvements compared to other programming languages are attributed to two factors. First, LLMs are proficient in Python, as evidenced by NAIVE attaining the highest average valid rate of 68.5% compared to others. Second, Python’s dynamic nature makes LSPAI difficult to detect errors before code execution, preventing LSP from identifying issues early and thereby degrading performance. Consequently, LSPAI decreased the valid rate for BAK and C4AI projects, indicating that LSP-guided diagnosis may not fully capture potential errors in Python tasks. Nonetheless, LSPAI successfully increased line coverage by generating unit tests in some cases. This shows that collected dependency information by LSPAI generated effective unit tests that cover more edge cases, resulting in more reliable test cases with higher coverage.

4.3 Breakdown of LSPAI

To evaluate the practical applicability of LSPAI in real-world use cases, we measured its time and token consumption for unit test generation. We put the averaged figure classified by programming languages. The statistics were collected using a maximum of five fixing attempts per focal method, with LSPAI accessing an LLM via API requests.

As shown in Table 4, LSPAI takes approximately 51 seconds and consumes 3,470 tokens on average to generate and refine unit tests. For Golang projects, the utility of LSPAI is particularly impressive, as it requires an acceptable range of resources (46 seconds and 4,384 tokens on average) while delivering an 8× improvement in line coverage. This demonstrates the efficiency and effectiveness of LSPAI in languages like Golang, where LLMs are less proficient, making automated unit test generation particularly effective. In contrast, LSPAI is less ideal for Python projects. While time and token usage remains within acceptable ranges, the modest improvement of approximately 6% in line coverage does not justify the resource expenditure of 1 minute and approximately 2,000 tokens per focal

Table 4: Time and Token Usage by LSPAI.

	Time(milliseconds)				Token	
	Retrieval	Diagnosis	GEN	FIX	GEN	FIX
Python	34363	9910	15745	1111	1639	103
Java	13065	3761	13728	19339	1316	2966
Golang	7941	3607	12742	24456	1229	3155
Averaged	18457	5759	14072	14968	1395	2075

method unless additional performance improvements in test quality can be ensured. This makes LSPAI less practical for Python projects compared to its strong performance in other languages. Overall, LSPAI demonstrates efficient and scalable performance across a variety of programming languages, but it excels most in contexts where LLMs traditionally underperform, such as with Golang. The results highlight LSPAI’s potential to bridge performance gaps in automated unit test generation for less-supported languages while maintaining acceptable resource usage for real-world applications.

5 Lessons Learned

In this section, we introduce some lessons learned during building tool for multi-language unit test generation and applying the tool to development environment in Tencent Ltd.

Alignment between Research and Industrial Needs. We identified a significant gap between academic research and industrial requirements in the domain of unit test generation. While academic efforts predominantly aim for high code coverage, they often overlook practical aspects essential for real-world usage. From our industry practice, developers urgently need a lightweight tool that can generate unit tests without whole-project compilation. Additionally, academic research has largely focused on specific programming languages such as Python and Java, leaving a gap in support for other languages. For example, Golang is widely adopted in many industries, yet few academic studies have addressed unit test generation for Golang. LSPAI was developed to bridge these gaps. Although it may not achieve the same level of code coverage as established academic tools like EvoSuite [10], LSPAI is popular among industrial developers.

Varying Language Proficiency of LLMs. Our experimental analysis revealed that LLMs exhibit varying levels of proficiency across different programming languages, directly impacting the effectiveness of LSPAI. Specifically, LLMs frequently make errors when generating Golang code, which affects the quality of the generated unit tests. These findings highlight the necessity for LLM-integrated tools to adapt their strategies to the specific demands and complexities of each programming language, thereby maximizing their utility and effectiveness.

Demands for Better Integration Methods. This work opens several promising avenues for future research in multi-language unit test generation through the LSP. Currently, LSPAI employs prompt engineering, a low-cost but limited method for harnessing the full potential of LLMs. A more sophisticated integration method is needed to build a retrieval system that can fully exploit the capabilities of LLMs. Besides, the integration of additional language

server functionalities, such as code intelligence and code action recommendations, could further enhance the accuracy and reliability of generated unit tests. In our industrial practice, we found that developers often require a more comprehensive approach to make the generated unit tests more reliable and useful.

6 Conclusion

We introduced LSPAI, a practical real-time unit test generation tool that leverages LLMs and integrates static analysis through the LSP to support multi-language codebases. LSPAI addresses the critical gap in existing research by enabling seamless unit test generation across diverse programming languages without the need for project-wide compilation, thereby facilitating concurrent test creation alongside code development. Implemented as an IDE plugin, LSPAI simplifies adoption for developers by requiring only the installation of appropriate language analysis tools. Our comprehensive evaluation of Java, Python, and Golang projects demonstrated that LSPAI consistently enhances both line coverage and valid rate compared to baseline approaches.

References

- [1] Juan Altmayer Pizzorno and Emery D Berger. Coverup: Coverage-guided llm-based test generation. *arXiv e-prints*, pages arXiv–2403, 2024.
- [2] Yiran Cheng, Lwin Khin Shar, Ting Zhang, Shouguo Yang, Chaopeng Dong, David Lo, Shichao Lv, Zhiqiang Shi, and Limin Sun. Llm-enhanced static analysis for precise identification of vulnerable oss versions. *arXiv preprint arXiv:2408.07321*, 2024.
- [3] Arghavan Moradi Dakhel, Amin Nikanjam, Vahid Majdinasab, Foutse Khomh, and Michel C. Desmarais. Effective test generation using pre-trained large language models and mutation testing. *Inf. Softw. Technol.*, 171(C), July 2024.
- [4] Amirhossein Deljouy, Roham Koohestani, Maliheh Izadi, and Andy Zaidman. Leveraging Large Language Models for Enhancing the Understandability of Generated Unit Tests. In *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*, pages 392–404, Los Alamitos, CA, USA, May 2025. IEEE Computer Society.
- [5] OpenAI et al. Gpt-4 technical report, 2024.
- [6] Apache Software Foundation. Apache commons cli, 2025. Accessed: 2025-01-15.
- [7] Apache Software Foundation. Apache commons csv, 2025. Accessed: 2025-01-15.
- [8] PSF (Python Software Foundation). Black, 2025. Accessed: 2025-01-15.
- [9] Steve Francia. Cobra. <https://github.com/spf13/cobra>, 2013. Accessed: 2025-01-15.
- [10] Gordon Fraser and Andrea 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*, pages 416–419, 2011.
- [11] Gordon Fraser and Andrea Arcuri. A large scale evaluation of automated unit test generation using evosuite. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 24(2):8, 2014.
- [12] Nat Friedman. Introducing github copilot: Your ai pair programmer, 2021. Accessed: 2024-12-20.
- [13] Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723*, 2020.
- [14] Gwhwan Go, Chijin Zhou, Quan Zhang, Xiazijian Zou, Heyuan Shi, and Yu Jiang. Towards more complete constraints for deep learning library testing via complementary set guided refinement. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2024*, page 1338–1350, New York, NY, USA, 2024. Association for Computing Machinery.
- [15] Nuno M Guerreiro, Duarte Alves, Jonas Waldendorf, Barry Haddow, Alexandra Birch, Pierre Colombo, and André FT Martins. Hallucinations in large multilingual translation models. *arXiv preprint arXiv:2303.16104*, 2023.
- [16] Caroline Lemieux, Jeevana Priya Inala, Shuvendu K Lahiri, and Siddhartha 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)*, pages 919–931. IEEE, 2023.
- [17] Ziyang Li, Saikat Dutta, and Mayur Naik. Llm-assisted static analysis for detecting security vulnerabilities. *arXiv preprint arXiv:2405.17238*, 2024.
- [18] Yun Lin, Jun Sun, Gordon Fraser, Ziheng Xiu, Ting Liu, and Jin Song Dong. Recovering fitness gradients for interprocedural boolean flags in search-based testing. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis*, pages 440–451, 2020.
- [19] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [20] Stephan Lukaszczuk and Gordon Fraser. Pynguin: Automated unit test generation for python. In *Proceedings of the ACM/IEEE 44th International Conference on Software Engineering: Companion Proceedings*, pages 168–172, 2022.
- [21] Microsoft. Best practices for using github copilot, 2024. Accessed: 2024-12-20.
- [22] Microsoft. Language server protocol, 2024. Accessed: 2024-12-24.
- [23] Microsoft. Syntax highlight guide, 2024. Accessed: 2024-12-20.
- [24] Chao Ni, Xiaoya Wang, Liushan Chen, Dehai Zhao, Zhengong Cai, Shaohua Wang, and Xiaohu Yang. Casmodatest: A cascaded and model-agnostic self-directed framework for unit test generation. *arXiv preprint arXiv:2406.15743*, 2024.
- [25] Carlos Pacheco and Michael D Ernst. Randoop: feedback-directed random testing for java. In *Companion to the 22nd ACM SIGPLAN conference on Object-oriented programming systems and applications companion*, pages 815–816, 2007.
- [26] Gabriel Ryan, Siddhartha Jain, Mingyue Shang, Shiqi Wang, Xiaofei Ma, Murali Krishna Ramanathan, and Baishakhi Ray. Code-aware prompting: A study of coverage-guided test generation in regression setting using llm. *Proceedings of the ACM on Software Engineering*, 1(FSE):951–971, 2024.
- [27] Arkadii Sapozhnikov, Mitchell Olsthorn, Annibale Panichella, Vladimir Kovalenko, and Pouria Derakhshanfar. Testspark: IntelliJ idea’s ultimate test generation companion. In *Proceedings of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion Proceedings*, pages 30–34, 2024.
- [28] Mohammed Latif Siddiq, Joanna Cecilia Da Silva Santos, Ridwanul Hasan Tanvir, Noshin Ulfat, Fahmid Al Rifat, and Vinicius Carvalho Lopes. Using large language models to generate junit tests: An empirical study. In *Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering, EASE '24*, page 313–322, New York, NY, USA, 2024. Association for Computing Machinery.
- [29] Sirupsen. Logrus, 2025. Accessed: 2025-01-15.
- [30] UncleCode. Crawl4ai, 2025. Accessed: 2025-01-15.
- [31] Han Wang, Han Hu, Chunyang Chen, and Burak Turhan. Chat-like asserts prediction with the support of large language model. *arXiv preprint arXiv:2407.21429*, 2024.
- [32] Siwei Wang, Xue Mao, Ziguang Cao, Yujun Gao, Qucheng Shen, and Chao Peng. Nxtunit: Automated unit test generation for go. In *Proceedings of the 27th International Conference on Evaluation and Assessment in Software Engineering, EASE '23*, page 176–179, New York, NY, USA, 2023. Association for Computing Machinery.
- [33] Zejun Wang, Kaibo Liu, Ge Li, and Zhi 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*, pages 1258–1268, 2024.
- [34] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [35] Zhuokui Xie, Yinghao Chen, Chen Zhi, Shuiguang Deng, and Jianwei Yin. Chatunitest: a chatgpt-based automated unit test generation tool. *arXiv preprint arXiv:2305.04764*, 2023.
- [36] Chen Yang, Junjie Chen, Bin Lin, Jianyi Zhou, and Ziqi Wang. Enhancing llm-based test generation for hard-to-cover branches via static analysis. *arXiv preprint arXiv:2404.04966*, 2024.
- [37] Zhichao Zhou, Yutian Tang, Yun Lin, and Jingzhu He. An llm-based readability measurement for unit tests’ context-aware inputs. *arXiv preprint arXiv:2407.21369*, 2024.

Received 2025-01-12; accepted 2025-03-26