# LSPRAG: LSP-Guided RAG for Language-Agnostic Real-Time Unit Test Generation

Gwihwan Go
Tsinghua University
Beijing, China
iejw1914@gmail.com

Quan Zhang*
East China Normal University
Shanghai, China
quanzh98@gmail.com

Chijin Zhou*
East China Normal University
Shanghai, China
tlock.chijin@gmail.com

Zhao Wei
Tencent
Beijing, China
zachwei@tencent.com

Yu Jiang
Tsinghua University
Beijing, China
jiangyu198964@126.com

## Abstract

Automated unit test generation is essential for robust software development, yet existing approaches struggle to generalize across multiple programming languages and operate within real-time development. While Large Language Models (LLMs) offer a promising solution, their ability to generate high coverage test code depends on prompting a concise context of the focal method. Current solutions, such as Retrieval-Augmented Generation, either rely on imprecise similarity-based searches or demand the creation of costly, language-specific static analysis pipelines. To address this gap, we present LSPRAG, a framework for concise-context retrieval tailored for real-time, language-agnostic unit test generation. LSPRAG leverages off-the-shelf Language Server Protocol (LSP) back-ends to supply LLMs with precise symbol definitions and references in real time. By reusing mature LSP servers, LSPRAG provides an LLM with language-aware context retrieval, requiring minimal per-language engineering effort. We evaluated LSPRAG on open-source projects spanning Java, Go, and Python. Compared to the best performance of baselines, LSPRAG increased line coverage by up to 174.55% for Golang, 213.31% for Java, and 31.57% for Python.

## CCS Concepts

• **Software and its engineering** → **Software testing and debugging**; • **Computing methodologies** → **Neural networks**; *Natural language processing*.

## Keywords

Unit Testing, Language Server Protocol, Retrieval Augmented Generation, Large Language Model

---

*Quan Zhang and Chijin Zhou are corresponding authors.

## 1 Introduction

Unit testing for modern software systems is crucial yet labor-intensive [4, 68]. To ensure reliability, test suites must achieve high coverage to expose subtle bugs and edge cases. It requires developers to have a deep understanding of the internal logic and dependencies of the codebase, including how the functions under test interact with their surrounding implementation. This challenge is even greater in today's software development, where enterprise codebases frequently span *multiple programming languages*, and developers expect *real-time* test generation as they code [1, 32, 41, 62].
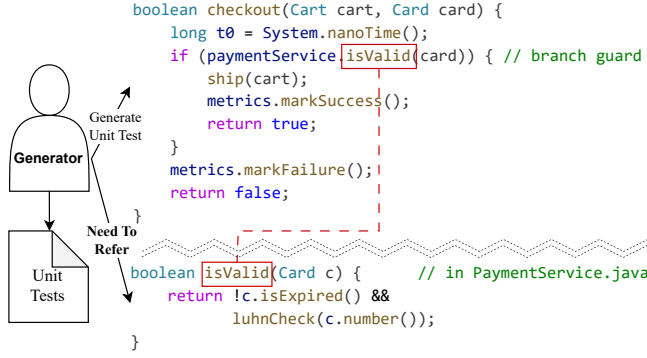
Recent advances in Large Language Models (LLMs) have invigorated research on real-time multi-language unit-test generation. Commercial systems such as GitHub Copilot [54] already offer one-click test generation, boosting test quality with simple heuristics—for example, giving higher weight to files a developer has recently opened [24]. At the same time, Repository-level Retrieval-Augmented Generation (RAG) techniques supply LLMs with additional code context obtained via textual similarity [38, 63], graph-based relationships [36, 37], or even web search [60, 66].

**Problem.** Although existing methods enhanced code generation practice, they are still limited in generating high coverage unit tests, because they *fail to retrieve relevant context precisely*. For example, in Figure 1, to cover the "true" branch of line 3 in `checkout` method, it is necessary to retrieve the definition of the guard-condition method `isValid`, which resides in another location. ( e.g., `PaymentService.java`). Existing approaches, however, struggle to handle this seemingly straightforward task effectively. GitHub Copilot is not able to pull dependent context across files, so the user must supply the relevant code manually [25].

Existing RAG approaches for code generation, unfortunately, also struggle to discover such a use-definition relationship precisely. This is because the RAG's embedding-and-similarity approach depends on superficial cues such as similar function names, variable names, and comments to infer relationships. These cues are often noisy and prone to variation across coding styles and conventions. Recent research [7, 36] recognizes this deficiency and augments retrieval with static-analysis-based searches; however, that solution demands substantial human effort and remains tied to specific programming languages, limiting their generalizability in today's

```java
boolean checkout(Cart cart, Card card) {
    long t0 = System.nanoTime();
    if (paymentService.isValid(card)) { // branch guard
        ship(cart);
        metrics.markSuccess();
        return true;
    }
    metrics.markFailure();
    return false;
}
```

```java
boolean isValid(Card c) {          // in PaymentService.java
    return !c.isExpired() &&
            luhnCheck(c.number());
}
```
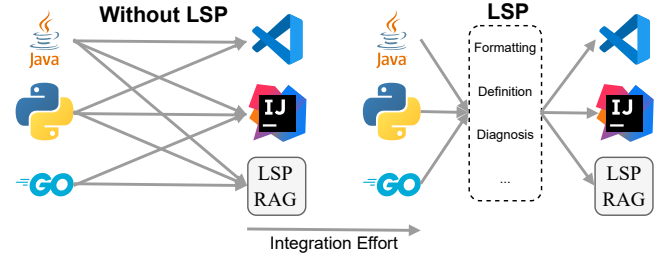
**Figure 1: Unit test generation scenario for focal method checkout. To cover the "true" block of the first branch of the focal method, the generator needs the definition of `isValid` from another file.**



**Figure 2: Integration effort before and after LSP.**

software development. This highlights the urgent need for a new approach for real-time and language-agnostic unit test generation.

**Insight.** This paper introduces LSPRAG, a concise context retrieval framework for multi-language unit test generation. Our key insight is that modern editors already ship with mature static analyzers exposed through the Language Server Protocol (LSP). By querying these analyzers on demand, LSPRAG obtains the precise location of context for every used symbol (e.g., function, method, etc) in the focal method.

**Challenges.** Despite having sufficient and precise symbol information, generating high-coverage unit tests remains non-trivial: there still exist challenges. (1) **The focal method includes excessive irrelevant context.** Achieving high test coverage requires precisely identifying the definitions of symbols that control branch conditions. However, in practice, a focal method often includes many symbols unrelated to the branch condition. For example, the nine-line `checkout` method in Figure 1 references eight external symbols (`Cart`, `Card`, `System`, `...`, and `markFailure`), each with potentially several lines of definition. Despite this, test coverage is only influenced by the `isValid` call. While retrieving the definition for `isValid` is essential, collecting every definition within the focal method introduces unnecessary, noisy context. This excess context complicates the generator's ability to locate the correct context for solving the branch condition, ultimately hindering the generation of high-quality unit tests.

(2) **Hard to guarantee valid tests in real-time.** Even when concise and relevant context is collected, simply providing LLMs with *more context does not guarantee* the generation of valid test code. As noted in prior research [35, 64, 67], large code context slices often introduce noise, which can obscure the original intent and reduce the LLM's ability to generate correct code. As a result, the generated test code is likely to contain syntax errors, which hinder high-coverage test generation rather than support it. Previous works [2, 23, 29, 47, 61] have mitigated this issue by executing the generated test code in advance and collecting error messages in a fixing loop. However, this approach is not practical for a real-time scenario, as it requires a significant amount of time to compile and

execute the generated test code. Furthermore, the codebase under test is often not compilable or executable in real-time.

We tackle these challenges from two dimensions. To address Challenge (1), we design a key token extraction strategy that identifies symbols essential for high-coverage test generation. Specifically, we employ a hybrid analysis approach that combines fine-grained lexical information from the LSP with structural information from the Abstract Syntax Tree (AST), effectively filtering out unnecessary context. To address Challenge (2), we design a compile-free self-repair mechanism that automatically fixes syntax errors in unit tests. Specifically, we utilize the LSP's diagnostic feature in a compilation-free loop and retrieve the necessary context for repairs. We continue this process iteratively until the errors disappear or the retry budget is exhausted, ensuring real-time performance.

We implemented LSPRAG as a Visual Studio Code Extension to streamline language-agnostic unit test generation and evaluated its performance on real-world Java, Python, and Golang projects. The results show that LSPRAG consistently improves the line coverage and the rate of valid tests generated, irrespective of the programming language or the underlying LLM employed. Our evaluation demonstrates that LSPRAG yields substantial improvements in unit test quality in terms of its line coverage and valid rate. Specifically, when compared to the best value among baselines, LSPRAG improved line coverage by a range of 91.4% to 174.55% for Golang projects, 27.79% to 213.31% for Java projects, and 16.87% to 31.57% for Python projects. Similarly, the valid rate of generated tests increased by up to 242.86% for Golang projects, 251.91% for Java projects, and 20.16% for Python projects.

In summary, this paper makes the following contributions:

- We identify a gap between academic research and industry practice: developers in large companies require high-quality test cases in real-time for different programming languages, yet concise context for LLM is difficult to obtain without compilation, execution, or heavyweight analysis.
- We designed and implemented LSPRAG, which generate language-agnostic high-coverage unit tests in real-time. The source code is available at https://thu-wingtecher.github.io/LSPRAG.
- We conducted a comprehensive evaluation of LSPRAG on real-world projects in three different programming languages. The results validate the ability of LSPRAG to consistently improve the performance of unit test generation in terms of both line coverage and the rate of valid tests.

## 2  Background

### 2.1  Language Server Protocol

*The Language Server Protocol* (LSP) was introduced by Microsoft [42] to address a significant challenge in software development: the need for each editor and IDE to have a unique implementation for every programming language's analysis features. As the number of editors and languages proliferated, this approach became unsustainable, forcing language authors to repeatedly rewrite the same analysis logic. At the same time, tool vendors struggled to keep pace with the emergence of new languages. LSP has resolved these problems by standardizing communication between a language server and any compliant editor or IDE.

**Before vs. After LSP** As illustrated in Fig 2, LSP significantly reduces implementation costs by defining a standardized communication protocol between language clients and language servers. The language client, typically an editor or IDE used by a developer (shown on the right in each sub-figure), communicates with a standalone language server process. This server, depicted on the left, handles parsing, static analysis, and other language-specific tasks. Before the adoption of LSP, adding rich language support required redundant effort. For instance, both Eclipse and VS Code needed separate Java plug-ins, with each language client implementing its pipeline for identical language features. LSP streamlines this process by establishing a standard protocol. Now, any compliant IDE can leverage existing language servers that also adhere to the LSP standard. This enables a single server to support multiple editors, allowing one editor to work with various languages by simply connecting to the appropriate servers.

**Typical LSP capabilities.** A standard LSP provides a suite of powerful features, including code completion, go-to-definition, find-all-references, and real-time diagnostics that display errors and warnings as you type. It also offers hover-to-view documentation, symbol renaming, code actions for quick fixes, and refactoring. Since these capabilities are delivered through a unified protocol, users can enjoy a consistent, IDE-grade experience across any LSP-compatible editor with minimal configuration.

### 2.2  Real-Time Unit Test Generation

In today's fast-paced software development, developers frequently need to generate unit tests while writing code, especially when working on incomplete or experimental features. The traditional approach of writing tests after the code is complete can slow down development, particularly when dealing with rapidly evolving projects. This is because some function implementations may be incomplete or stubbed, which can temporarily leave the project not compilable and complicate test generation. Real-time unit test generation addresses this by enabling developers to generate tests concurrently with their code, even for incomplete or partial function implementations. Modern tools like GitHub Copilot [54] and Amazon CodeWhisperer [51] leverage LLMs to create unit tests in seconds, helping developers stay in their flow. These are integrated directly into IDEs, providing test candidates within well under a minute, ensuring that the development process remains uninterrupted.
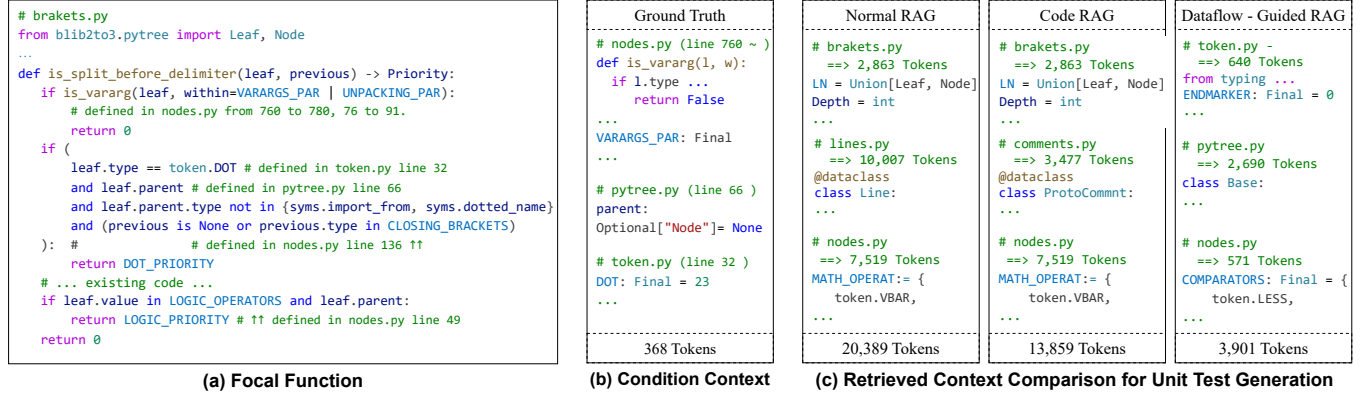
## 3  Motivating Example

Existing real-time code-generation techniques struggle to produce comprehensive unit tests, mainly because they cannot automatically retrieve the concise context that the focal method relies on. Closing this gap is essential because exercising every branch demands a precise understanding of the symbols that guard those branches. Without this understanding, it's difficult for LLM to understand branch conditions and fails to generate the specific inputs needed to exercise them. To ground the problem, we present a motivating example from the project Black [21], a widely used Python code formatter, and compare how existing research efforts attempt to retrieve necessary context.

Figure 3 (a) illustrates the simplified focal method of the Black project for which we want to generate unit tests. The method `is_split_before_delimiter` evaluates a tree leaf and its predecessor to assign a numeric priority that determines whether the formatter should split the line before that delimiter. It skips var-arg constructs, then gives higher priority to dots in attribute chains and logical operators, returning 0 when no split is warranted. To trigger diverse actions of this method, the generator must supply inputs that make each branch guard evaluate to true and false. For instance, to satisfy the first guard, we must understand exactly when the function `is_vararg` returns true and which value of leaf produces that outcome. To do this, we need to retrieve the definition of the function `is_vararg` from other files, which requires (1) identifying the correct file names and (2) isolating the precise code snippets that define the needed functions.

Context retrieval is a well-studied problem in the field of LLM, with many RAG techniques proposing advanced strategies. Existing RAG techniques may address real-time unit-test generation for the focal method `is_split_before_delimiter` in three main ways. (i) Normal RAG [30] employs similarity search with fixed-size chunking, designed for unstructured natural language. (ii) Code-aware RAG, such as CodeRAG [45] incorporates code-specific features by employing AST-based chunking and codebase indexing. (iii) Static-analysis-based RAG, such as DraCo [7], enhances code-aware RAG by leveraging traditional program analysis techniques, such as data-flow analysis. Although RAG techniques have shown encouraging results, they still face limitations in retrieving the necessary context for unit test generation.

These limitations are exemplified in Figure 3 (c). Both Normal RAG and Code RAG were able to identify only partial dependencies (e.g., `nodes.py`) while missing other essential contexts (e.g., `pytree.py` and `tokens.py`). This occurs because their retrieval is based on an embedding-and-similarity algorithm, which depends on textual information such as variable names, function names, and comments. These cues are often noisy and prone to variation across different coding styles and project conventions. Furthermore, they cannot isolate the relevant snippet and therefore feed the LLM large amounts of noise (e.g., 20,389 and 13,859 tokens) that hinders code generation. As the latest open-sourced technique for static-analysis-based RAG, DraCo addressed the above limitations by building an import graph to follow dependencies, yet notable drawbacks remain. First, it is still prone to including noise. Because it cannot recognize which part of the imported library is used, it retrieves every imported target, thereby including unnecessary context. For instance,

**Figure 3: Illustrative example of existing work's context retrieval in unit-test generation. (a) presents a simplified version of the focal method; (b) displays the ground-truth condition context and its location; and (c) compares the context slices and the number of tokens retrieved by three existing RAG techniques.**

it loads all constant definitions from tokens.py when only a single constant, token.DOT, is required. Second, the implementation cost is prohibitive. Although DraCo's retrieval accuracy is relatively high, constructing and maintaining these language-specific data-flow graphs requires substantial manual effort, limiting DraCo to codebases written in Python alone.

**Our approach.** Rather than hand-crafting language-specific data-flow graphs, we simply reuse the static-analysis features of LSP. Through the LSP, these services surface the exact, minimal code fragment for requested symbols, in almost any mainstream language, giving us precise context with minimal effort. The next section details how LSPRAG turns this observation into practice.

## 4 Basic Concepts

This section formalises the four LSP providers used by LSPRAG. We introduce each concept in a top-down hierarchy—from workspace to file, symbol, and token—so that their relationships are clear before provider interfaces are stated.

### 4.1 Workspace, File, Symbol, and Token

*Definition 4.1 (Workspace).* A *workspace*, denoted $\mathcal{W}$, is a finite set of *source files*: $\mathcal{W} = \{ f_1, \ldots, f_{|\mathcal{W}|} \}$.

*Definition 4.2 (File).* Each file $f \in \mathcal{W}$ is an ordered sequence of characters and contains a finite set of *symbols*, $\Sigma_f = \{\sigma_1, \ldots, \sigma_{|\Sigma_f|}\}$.

*Definition 4.3 (Symbol).* A *symbol* $\sigma$ is a named program entity, such as a class, function, or constant, represented by the tuple

$$\sigma = (name, kind, loc, children).$$

Here, *name* is the identifier appearing in the source code; *kind* is an enumerated tag (e.g., Class, Function, or Constant); *loc* is the inclusive character location $(f, o_{\text{start}}, o_{\text{end}})$ of the symbol within its file $f$, start offset $o_{\text{start}}$, and end offset $o_{\text{end}}$; and *children* $\subseteq \Sigma_f$ is the set of nested symbols whose spans lie strictly inside *loc* (e.g., a method symbol is one of the children from its class symbol.)

*Definition 4.4 (Token).* A *token* $\tau$ is the smallest lexical unit recognised through LSP. We describe $\tau$ as the tuple

$$\tau = (loc, lex, tok\_kind), \qquad f \in \mathcal{W}.$$

The $f$ denotes the file containing the token; *loc* is the character location $(f, o_{\text{start}}, o_{\text{end}})$; *lex* is the raw lexeme text; and *tok_kind* indicates the lexical category (e.g., identifier, keyword, …).

### 4.2 Provider Interfaces

In addition to the above concepts, we specify the four code-context providers implemented on top of LSP that LSPRAG invokes.

*Definition 4.5 (Symbol Provider).* The *Symbol Provider*, denoted *SYM*, is a function that, given a file $f \in \mathcal{W}$, returns the set of all symbols in $f$ that are not nested within another symbol.

$$SYM(f) = \Sigma_f \setminus \bigcup_{\sigma' \in \Sigma_f} \sigma'.children$$

where $\Sigma_f$ is the set of all symbols in file $f$.

*Definition 4.6 (Token Provider).* The *token provider*, denoted *TOK*, is a function that, given a location $loc_q \in \mathcal{W}$, returns an ordered sequence of all tokens whose spans are contained within $loc_q$. Let $T_f$ be the complete, ordered sequence of all tokens in file $f$. The function is defined as:

$$TOK(loc_q) = \langle \tau \in T_f \mid \tau.loc \subseteq loc_q \rangle$$

The output is a sequence, not a set, preserving the original order of the tokens as they appear in the file.

*Definition 4.7 (Definition Provider).* A *Definition Provider*, denoted *DEF*, is a function that maps a given token $\tau$ to the set of locations $L_\tau$ of its corresponding definition symbols. This can be expressed as:

$$DEF(\tau) = L_\tau \triangleq \begin{cases} \{\sigma.loc \mid \sigma \text{ is a definition for } \tau\} & \text{if definitions exist} \\ \emptyset & \text{otherwise} \end{cases}$$

Subsequent to the definition retrieval, resolving a location $l$ to a symbol $\sigma$ involves finding the symbol within the file's symbol set from $SYM(l.f)$ that satisfies:

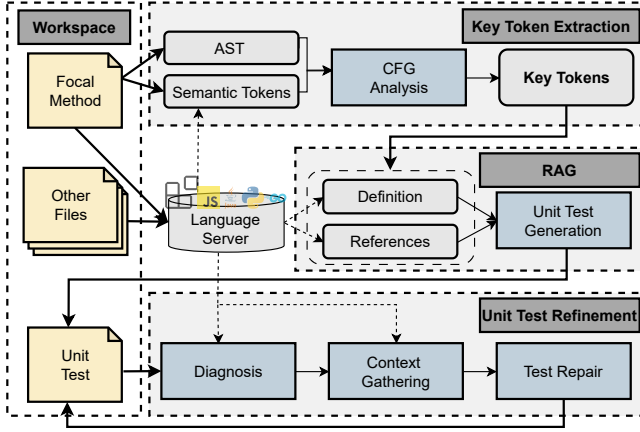$$\sigma \in SYM(l.f) \land \sigma.loc = l$$

**Figure 4: Overall workflow of LSPRAG.**

Note that the file containing the definition, $l.f$, may reside outside of the current workspace $\mathcal{W}$, for instance, within an external library or a standard system header.

*Definition 4.8 (Reference Provider).* The *reference provider*, denoted *REF*, performs the inverse operation of the definition provider. Given a token $\tau$, it finds all tokens in the workspace that refer to the same definition. Let $\sigma_{def}$ indicates the definition symbol for the token $\tau$. The function is defined as:

$$REF(\tau) = \begin{cases} \{\tau'.loc \mid \tau' \in \bigcup_{f \in W} T_f \\ \qquad \wedge DEF(\tau') = \sigma_{def}\} & \text{if } \sigma_{def} \neq \texttt{null} \\ \emptyset & \text{if } \sigma_{def} = \texttt{null} \end{cases}$$

The output is a set of locations, each corresponding to a token that references the same definition symbol.

## 5 Design of LSPRAG

This section presents the design of LSPRAG, an enhanced RAG framework for multi-language real-time unit test generation. As illustrated in Figure 4, LSPRAG consists of three core modules that operate on a developer's workspace $\mathcal{W}$. We assume that a developer is working on a focal method within a file $f \in \mathcal{W}$, which may have dependencies on various symbols defined in other files.

Upon a developer's request to generate a unit test, LSPRAG begins with **Key Token Extraction** module. This module is responsible for identifying a set of key tokens within the focal method. These tokens are pivotal in guiding the generation of high-coverage test cases as they typically govern the method's control flow and its interactions with external components. Subsequently, these extracted tokens are passed to the **RAG** Module. This module leverages *DEF* and *REF* providers to gather contextual information, including definitions and usages of the key tokens. This context is then used to construct a detailed prompt that is fed to an LLM. Finally, to address potential syntactic errors in the generated test code, the **Unit Test Refinement** module inspects the test code without compilation. If any potential errors are detected, it gathers the necessary context for remediation and feeds this information back to the LLM to correct the code.

## 5.1 Key Token Extraction

A primary challenge in context retrieval is that a focal method often contains many tokens irrelevant to its core logic. Including context of these tokens can degrade the quality of the generated tests by retrieving non-essential context. To address this, we introduce the concept of a *key token*, which is a token $\tau$ whose semantics are essential for constructing high-coverage tests. Specifically, we define key tokens as those that represent control-flow decisions or external dependencies. To identify these tokens, LSPRAG performs a hybrid analysis that combines fine-grained lexical information with a structural understanding, as illustrated in Figure 5.

*5.1.1 Lexical Information from LSP.* LSPRAG begins by invoking a token provider, *TOK* (Definition 4.6), to retrieve a sequence of all tokens within the scope of the focal method's location, $loc_m$.
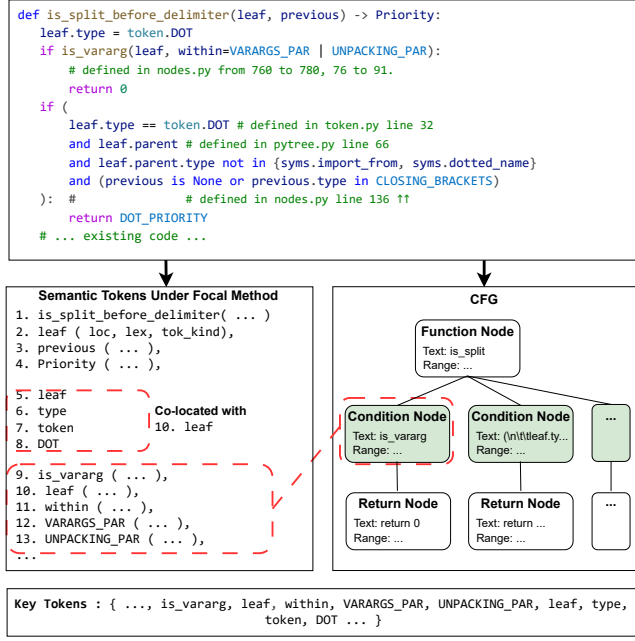
$$\langle \tau_1, \tau_2, \ldots, \tau_n \rangle = TOK(loc_m)$$

Each token $\tau$ in this sequence includes its location, lexical name, and a semantic role, such as parameter or identifier. This information is necessary for the subsequent context retrieval phase, because the context provider on top of LSP accepts input as a specific location of a token. Additionally, based on semantic role, one can filter out unnecessary tokens, such as constants. While this lexical information is useful, it is insufficient on its own because it lacks structural context. For example, tokens like "if" and "def" are both classified simply as "Keyword" by *TOK*, despite their different structural implications for control flow. The LSP-oriented solution cannot resolve this limitation, since it is designed for editor-centric features, not compiler-level control-flow analysis[1]. Consequently, relying solely on LSP features is insufficient for identifying the key tokens needed for high-coverage test generation.

*5.1.2 Structural Analysis via AST.* To overcome the limitations of the LSP, LSPRAG leverages the Abstract Syntax Tree (AST). The AST provides a language-agnostic structural view that is essential for understanding program hierarchy and can be easily implemented in a language-agnostic way, which is crucial for multi-language support. The effectiveness of this technique has been demonstrated in prior RAG research [7, 50]. LSPRAG reuses widely-used Tree-sitter [40] to construct a AST for focal method. The critical step in this process is to bridge the structural AST view with the lexical LSP tokens. To achieve this, LSPRAG preserves the precise source code location (i.e., line and column numbers) of every node during AST construction. This allows for a precise mapping between an AST node and its corresponding set of tokens if the node's location span contains the tokens' locations.

*5.1.3 Key Token Identification.* With the combined data, the tool constructs a lightweight, language-agnostic Control-Flow Graph (CFG) for the focal method. The CFG enables path-sensitive analysis to identify the final set of key tokens, $T_{\text{key}}$. At a high level, a token $\tau$ from the focal method is added to $T_{\text{key}}$ if it participates in or affects the change of a conditional expression that determines a branch in the CFG. Specifically, we first identify tokens $\tau'$ that participate in a conditional expression. Next, we filter out the token sequence $TOK(loc_m)$ that does not need a context search by excluding tokens with semantic roles such as Keyword, Identifier,

---

[1]https://github.com/microsoft/language-server-protocol/issues/1675

```
def is_split_before_delimiter(leaf, previous) -> Priority:
    leaf.type = token.DOT
    if is_vararg(leaf, within=VARARGS_PAR | UNPACKING_PAR):
        # defined in nodes.py from 760 to 780, 76 to 91.
        return 0
    if (
        leaf.type == token.DOT # defined in token.py line 32
        and leaf.parent # defined in pytree.py line 66
        and leaf.parent.type not in {syms.import_from, syms.dotted_name}
        and (previous is None or previous.type in CLOSING_BRACKETS)
    ): #           # defined in nodes.py line 136 ↑↑
        return DOT_PRIORITY
    # ... existing code ...
```

**Semantic Tokens Under Focal Method**

```
1. is_split_before_delimiter( ... )
2. leaf ( loc, lex, tok_kind),
3. previous ( ... ),
4. Priority ( ... ),

5. leaf
6. type            Co-located with
7. token           10. leaf
8. DOT

9. is_vararg ( ... ),
10. leaf ( ... ),
11. within ( ... ),
12. VARARGS_PAR ( ... ),
13. UNPACKING_PAR ( ... ),
...
```

**CFG**

**Function Node**
Text: is_split
Range: ...

**Condition Node**
Text: is_vararg
Range: ...

**Condition Node**
Text: (\n\t!leaf.ty...
Range: ...

...

**Return Node**
Text: return 0
Range: ...

**Return Node**
Text: return ...
Range: ...

...

**Key Tokens** : { ..., is_vararg, leaf, within, VARARGS_PAR, UNPACKING_PAR, leaf, type, token, DOT ... }

**Figure 5: This figure illustrates the key token selection process. The red box highlights how key tokens are identified. Using the lexical information and CFG, LspRag identifies those tokens involved in or influencing branch decisions.**

Literal, Comment, String, and Regex. In general, these tokens are either constants or syntax keywords that do not require a context search. Then, we extract tokens from $TOK(loc_m)$ that appear on the same line as the token $\tau'$ and add them to $T_{key}$. The resulting set $T_{key}$ represents the tokens most likely to affect the focal method's execution path. It is then passed to the subsequent RAG module.

## 5.2 RAG

This module is designed for precise context retrieval for the set of key tokens, $T_{key}$, identified in the previous stage. To achieve this, LspRag systematically queries the LSP's *DEF* (Definition 4.7) and *REF* (Definition 4.8) for each token in $T_{key}$ and assembles the retrieved information into a coherent prompt.

*5.2.1 Context Retrieval.* To retrieve a token's definition, LspRag invokes the *DEF* to find the location where the token is defined, which we denote as $loc_{def}$. By default, *DEF* resolves symbols across the entire project space, including standard and third-party libraries. To maintain contextual relevance, LspRag filters these results, considering only definitions within the current workspace $\mathcal{W}$. This prevents the inclusion of unnecessary definitions from standard libraries (e.g., String, Integer). This filtering is expressed as:

$$\text{process } loc_{def} \text{ only if } loc_{def}.f \in \mathcal{W}, \text{ where } loc_{def} = DEF(\tau.loc)$$

If a valid definition location is found, LspRag first locates $\sigma_{def}$ at $loc_{def}$ and then extracts the source code from $\sigma_{def}.loc$.

Next, to find usage examples, LspRag queries the *REF* to find all use-sites of the token's symbol within the workspace $\mathcal{W}$. We denote $L_{refs}$ as the queried result of $REF(\tau)$. A raw reference location

$loc'_r \in L_{refs}$ (e.g., a variable in an expression) is often insufficient, as it only points to the token at other file without its surrounding context. To create a meaningful example, LspRag enriches each reference. For each location $loc'_r \in L_{refs}$, the tool identifies the smallest symbol $\sigma_{enclosing}$ that fully contains the reference location ($loc'_r \subseteq \sigma_{enclosing}.loc$). This is achieved by retrieving all top-level symbols in the reference's file $loc'_r.f$ through *SYM* and then recursively searching their children. The source code of $\sigma_{enclosing}$ is then extracted, transforming a simple location into a complete usage example (e.g., the entire function that uses the token $\tau$), which demonstrates how the symbol is used in other locations.

*5.2.2 Prompt Construction.* The final step synthesizes the retrieved information into a structured prompt for the LLM. This prompt is designed to guide the LLM in generating comprehensive unit tests and consists of three parts: (1) The full source code of the focal method to be tested. (2) The definitions and references of the key tokens retrieved previously. (3) A lightweight unit-test template with necessary import statements and class or function structures inferred from the focal method's file. This structured prompt provides the LLM with focused and relevant information, minimizing noise that can degrade code generation quality and empowering it to produce accurate and thorough unit tests. Due to page limitations, we provide the full prompt format in our artifact, which is open-sourced. The prompt is designed using CoT [59] with a one-shot setting [5].

## 5.3 Unit Test Refinement

While LLMs excel at code generation, their output can be imperfect, containing syntactic errors that cause compilation failures. To address this, LspRag employs an automated refinement module that detects and repairs them. This process operates in two phases: (1) real-time error detection and (2) context-aware error repair.

*5.3.1 Real-Time Error Detection.* LspRag leverages the diagnostic capabilities of the LSP to perform immediate analysis on the generated code. In detail, when code is modified within a file $f \in W$, LspRag notifies the corresponding language server. The server, in turn, analyzes the file and returns a set of diagnostics. Each diagnostic is a tuple $(m, loc)$, where $m$ is a human-readable error message (e.g., "undefined variable 'x'") and $loc$ is the location of the code that triggers the error. This mechanism provides instantaneous feedback without requiring a project compilation or code execution, which is critical for real-time context.

**Table 1: Error Categories Grouped by Necessary Context**

| Error Category | Frequency | Context Needed for Fix |
|---|---|---|
| Redeclaration/Duplicate Definition | 28 300 | Workspace-level context |
| Import/Module Resolution Error | 13 517 | (e.g., project structure) |
| Member Access/Usage Error | 13 387 | Symbol-level context |
| Type Mismatch/Compatibility Error | 4388 | (e.g., context of a specific symbol) |
| Constructor Call Error | 1670 | |
| Syntax Error | 5467 | No external context required |
| Unhandled Exception | 179 | |
| **Total** | **66908** | |

*5.3.2 Context-Aware Error Repair.* Upon receiving a set of non-empty diagnostics, LSPRAG initiates the repair phase. The core challenge in automated code repair is providing the LLM with sufficient and relevant context to understand and fix the error.

To determine the optimal context for different error types, We constructed the dataset (Table 1) ourselves by collecting 70k errors across Java, Python, and Golang, since no dataset exists for LLM-generated code paired with error messages. These errors were generated by prompting diverse LLMs to produce test cases for different projects (including those used in evaluation), and then diagnosing them via the LSP-based workflow. After collection, we manually analyzed the errors, identifying their root causes and the contextual information helpful in fixing them. This information is general rather than tied to specific projects. For instance, "imported and not used" consistently indicates import-related issues, while "redeclared in this block" points to symbol redeclaration errors. Based on the analysis result, LSPRAG gathers the corresponding necessary context by invoking the providers implemented on top of LSP. For a given diagnostic $(m, loc_{err})$, LSPRAG collects context as follows:

- **Symbol-Level Context:** For errors related to a specific symbol (e.g., undefined variables, incorrect method calls), LSPRAG retrieves its definition and usage examples. It invokes *DEF* to fetch the canonical definition and *REF* to find usage examples for the tokens at $loc_{err}$, providing the LLM with both the ground-truth specification and concrete examples of correct usage.
- **Workspace-Level Context:** For errors involving project structure (e.g., missing imports), a broader view of the project is necessary. LSPRAG provides this by supplying: (1) the workspace's file structure and (2) a list of top-level symbols (e.g., classes, functions) in the file where the error occurred $(loc_{err}.f)$.

Finally, LSPRAG constructs a prompt containing the erroneous code, the diagnostic message $m$, and the retrieved context. This process is iterative: after the LLM proposes a fix, LSPRAG re-triggers the diagnostic mechanism. This loop continues until all diagnostics are resolved or a predefined number of attempts is reached.

## 6  Implementation

We implemented a prototype of LSPRAG in TypeScript, comprising approximately 22k lines of code out of 8k lines of unit test code. The implementation was guided by two primary goals: developer usability and language generalizability.

To enhance usability, LSPRAG was implemented as a VS Code extension, which integrates seamlessly into the developer's existing workflow and reduces the barriers to adoption. The tool leverages the VS Code Extension API, which provides comprehensive capabilities for user interface interactions, editor functionality, and command handling. To achieve language generalizability, we designed LSPRAG to be language-agnostic. It leverages the LSP to consume semantic information from any compliant language server and utilizes Tree-sitter [40] for robust source code parsing into an AST. Consequently, LSPRAG is compatible with any programming language supported by an LSP server and Tree-sitter. We have validated LSPRAG's functionality with several major language servers, including Pylance for Python [43], the Oracle Extension Pack for Java [46], and the official Go extension [27].

LSPRAG currently offers two main features for unit test generation. The first is method-level test generation, which automatically generates a unit test class for a given focal method. Such a class contains multiple test functions, each designed for a unique execution branch within the focal method. The second feature, branch-targeted test generation, allows developers to select a specific conditional branch within a focal method and generates a set of test cases specifically designed to exercise the logic of that branch.

## 7  Evaluation

In this section, we comprehensively evaluate LSPRAG's performance on real-world projects across different programming languages. Our evaluation investigates the following questions:

- **RQ1**(§7.2): Can LSPRAG generate *higher coverage* unit tests than other baselines?
- **RQ2**(§7.3): Does LSPRAG maintain low latency and efficient token usage during generation?
- **RQ3**(§7.4): How much does each pipeline component of LSPRAG contribute to test coverage?

### 7.1  Experiment Setup

**Benchmarks.** Our evaluation is conducted on six open-source projects written in Python, Java, and Golang. Python and Java were selected due to their common use in prior unit test generation research. Golang was chosen as it is widely adopted in industry but less examined in academic studies, which we anticipated would present different challenges for LLMs. As shown in Table 2, we selected two projects per language, all of which are frequently used as benchmarks in the unit test generation domain [29, 39, 57, 58, 61].

**Table 2: Dataset Statistics**

| Project | Domain | Version | Language |
|---------|--------|---------|----------|
| Black [21] | Code formatter | 8dc9127 | Python |
| Tornado [14] | Network/Web framework | 81d36df1 | Python |
| Commons-CLI [19] | Command line interface parser | 266ab84a | Java |
| Commons-CSV [20] | Library for processing CSV files | ca3a95c3 | Java |
| Cobra [22] | Framework for creating Go CLI | ceb39ab | Golang |
| Logrus [53] | Structured logging library | d1e6332 | Golang |

**Baseline Selection.** To evaluate the performance of LSPRAG, we selected four baselines, comprising three RAG techniques and one prompt engineering method.

- RAG: A normal RAG with OpenAI embeddings and FAISS index via LangChain. It retrieves the top 3 most similar snippets based on cosine vector similarity.
- CODEQA [50]: An open-source RAG pipeline designed for code generation. Its approach, which involves indexing the codebase and parsing code structure with an AST, is analogous to that used in modern AI-assisted IDEs like Cursor [11].
- DRACO [7]: A RAG pipeline integrated with program analysis. Although it is specific to Python, this tool's use of static analysis-guided retrieval plays a comparable baseline for evaluating the context retrieval capabilities of LSPRAG.
- SYMPROMPT [48]: A prompt engineering technique for unit test generation. It proposes a novel prompt format that translates

**Table 3: Comparison of line coverage and valid rate across all baselines. Our approach, LspRag, consistently outperforms all baselines, achieving higher scores across diverse projects and models. Missing values (–) indicate that the baseline does not support test generation for those projects.**

| Project | Model | Coverage | | | | | Valid Rate | | | | |
|---------|-------|----------|------|-----------|--------|-------|----------|------|-----------|--------|-------|
| | | CodeQA | RAG | SymPrompt | LspRag | DraCo | CodeQA | RAG | SymPrompt | LspRag | DraCo |
| Black | GPT4o-m | 34.92% | 26.29% | 22.09% | **41.91%** | 32.58% | 55.57% | 37.26% | 66.35% | **79.73%** | 54.52% |
| | GPT4o | 40.47% | 29.94% | 23.99% | **48.10%** | 35.33% | 79.46% | 77.53% | 59.87% | **85.15%** | 73.31% |
| | DS-v3 | 46.03% | 38.98% | 34.18% | **56.89%** | 37.64% | 59.26% | 72.24% | 75.18% | **88.03%** | 71.04% |
| Tornado | GPT4o-m | 29.04% | 49.44% | 31.44% | **57.78%** | 46.54% | 71.26% | 64.30% | 77.01% | **81.84%** | 75.82% |
| | GPT4o | 29.97% | 46.29% | 34.35% | **60.90%** | 41.09% | 80.14% | 83.95% | 72.82% | 84.53% | **85.03%** |
| | DS-v3 | 38.45% | 54.05% | 50.78% | **64.62%** | 48.99% | 85.93% | 81.84% | 85.30% | **91.13%** | 90.48% |
| Commons CLI | GPT4o-m | 10.63% | 5.02% | 2.76% | **33.29%** | – | 12.40% | 8.18% | 7.05% | **45.08%** | – |
| | GPT4o | 9.60% | 3.21% | 3.00% | **34.69%** | – | 8.23% | 7.27% | 7.20% | **48.12%** | – |
| | DS-v3 | 20.72% | 17.68% | 5.62% | **37.72%** | – | 13.29% | 14.55% | 9.21% | **60.52%** | – |
| Commons CSV | GPT4o-m | 40.97% | 24.93% | 18.50% | **80.51%** | – | 23.64% | 13.20% | 6.28% | **82.85%** | – |
| | GPT4o | 44.85% | 44.84% | 25.28% | **78.33%** | – | 20.68% | 26.53% | 14.41% | **90.90%** | – |
| | DS-v3 | 65.11% | 44.68% | 35.07% | **83.20%** | – | 43.27% | 32.24% | 29.82% | **90.95%** | – |
| Cobra | GPT4o-m | 7.11% | 12.03% | 3.39% | **23.03%** | – | 6.01% | 9.50% | 1.23% | **23.88%** | – |
| | GPT4o | 10.05% | 7.57% | 0.22% | **27.60%** | – | 9.70% | 8.12% | 0.81% | **33.27%** | – |
| | DS-v3 | 15.50% | 13.01% | 8.57% | **37.22%** | – | 10.30% | 10.69% | 2.78% | **34.65%** | – |
| Logrus | GPT4o-m | 5.52% | 11.14% | 0.23% | **23.71%** | – | 14.32% | 20.83% | 0.83% | **34.02%** | – |
| | GPT4o | 5.61% | 13.09% | 0.23% | **27.75%** | – | 14.17% | 26.52% | 0.83% | **32.02%** | – |
| | DS-v3 | 11.34% | 11.00% | 5.43% | **21.81%** | – | 13.33% | 15.83% | 7.5 % | **33.11%** | – |

CFG paths into descriptive comments, thereby improving the readability and structure of the generated tests. We replicated this prompt format using our own CFG.

**Model Selection and Environment.** Our evaluation utilizes three LLMs to assess LspRag's performance across different architectures and scales. We selected GPT-4o [18] and its smaller variant, GPT-4o-mini (GPT4o-m), to represent the transformer architecture. To test against a different model design, we also included Deepseek-V3 [34] (DS-v3), which is based on a Mixture-of-Experts architecture. The experiments were conducted on a server running Ubuntu 22.04 LTS, equipped with a 128-core AMD EPYC 7763 CPU and an NVIDIA V100 GPU featuring 32GB of memory. For all language models, we used the default generation parameters, including temperature, to ensure a consistent baseline.

## 7.2 Generation of High Coverage Unit Tests

To answer RQ1, we evaluated LspRag's ability to generate high coverage unit tests against selected baselines. We measured performance using two key metrics: line coverage to quantify test effectiveness, and valid rate to assess the proportion of syntactically correct, usable tests. We define the valid rate as the ratio of grammatically correct test cases to the total number of unit test generation attempts. Since invalid test cases are not compilable or executable, they cannot produce any coverage.

**Evaluation Setup.** For each project listed in Table 2, we generated unit tests for methods exceeding 10 lines of code, as these typically contain complex logic that is difficult to fully cover. This

resulted in 299 focal methods for Black, 521 for Tornado, 101 for Cobra, 24 for Logrus, 43 for Commons-CLI, and 145 for Commons-CSV. To ensure reliable results, we repeated each experiment five times and report the average values. We configured LspRag to perform a maximum of five self-correction iterations.

**Statistical Analysis.** To verify the significance of our findings, we performed statistical tests on the results. We computed p-values to assess the significance of the improvements and reported effect sizes to quantify their magnitude, following the guidelines from Arcuri et al [3]. The null hypothesis ($H_0$) posits that there is no statistically significant difference in line coverage between LspRag and the baselines. Conversely, the alternative hypothesis ($H_1$) states that LspRag achieves a statistically significant improvement.

**Results.** The overall results, presented in Table 3, show that LspRag consistently improves line coverage across all projects, regardless of the programming language or the underlying LLM. Our statistical analysis confirms that these improvements are significant, with all p-values being less than 0.05, thus supporting the $H_1$. This finding highlights the effectiveness of LspRag's context retrieval and code correction mechanisms. The following analysis examines the specific design choices in the baselines that may contribute to this performance gap.

While SymPrompt is designed to guide an LLM in generating tests for specific conditions specified within comments, we found its effectiveness in our experiments to be limited by the absence of code dependency information. We observed that its complex prompt structure, when not grounded with the necessary context,

complicated the LLM's generation process. This suggests that the approach is robust but requires sufficient contextual input to steer the LLM effectively, and without it, it can struggle to match the performance of simpler methods.

The RAG-based approaches (RAG, CODEQA, and DRACO) showed varied performance, highlighting different trade-offs in context retrieval. For instance, DRACO is a sophisticated tool employing dedicated static analysis for Python. While effective, its reliance on import statements for dependency analysis means it can sometimes miss crucial intra-file context, such as helper functions or constants defined in the same file. On the other hand, CODEQA and RAG are designed to retrieve full code snippets of functions or classes, including intra-file context. This approach, however, sometimes led to an excess of retrieved context, particularly in larger projects, such as Tornado. This frequently caused the context to exceed the LLM's token limitation, resulting in the loss of potentially important information. Interestingly, our simpler baseline RAG, which retrieves a maximum of three code snippets, often performed better than the more complex RAG tools, suggesting that a more constrained and focused context can be more effective than a larger, unverified one.

LSPRAG's strong performance stems from two primary mechanisms. First, for higher coverage, LSPRAG leverages LSP features like "go to definition" and "find references" to construct a precise semantic dependency graph. This provides the LLM with a complete yet concise context, avoiding the noise of generic RAG and the intra-file blind spots of other methods. Second, to achieve a higher valid rate, LSPRAG utilizes real-time diagnostics on top of LSP to detect and correct syntactic errors in the generated code iteratively. This self-correction loop consistently improves the valid rate of generated test cases, as evidenced by the consistently high valid rates in Table 3. The specific contribution of these LSP-driven features is quantified in our ablation study (§7.4).

The improvement in the valid rate for the Python projects (Black and Tornado) was less obvious than for the Java and Golang projects. This is an expected outcome, as Python is a dynamically typed language, which limits the language server's ability to catch all potential errors statically before execution. Conversely, for projects written in statically-typed languages like Java and Go, LSPRAG consistently delivered significant improvements in validity rates because the language server could identify a broader range of potential errors.

We did not include Copilot in our experiments because it heavily relies on human interaction during the unit test generation. In detail, it requires manual clicks to trigger its functionality, which prevents large-scale automated experimentation as it does not provide automated APIs for programmatic code generation [9, 10]. Furthermore, it assumes that users will perform context collection and error correction during interaction, introducing human factors that make fair and consistent evaluation difficult.

## 7.3 Latency and Cost

For a code generation tool to be practical in an interactive development workflow, it must operate with acceptable latency and cost. To answer RQ2, we assess whether LSPRAG meets this real-time requirement by measuring its performance overhead in terms of latency and cost. The evaluation was conducted on the projects from
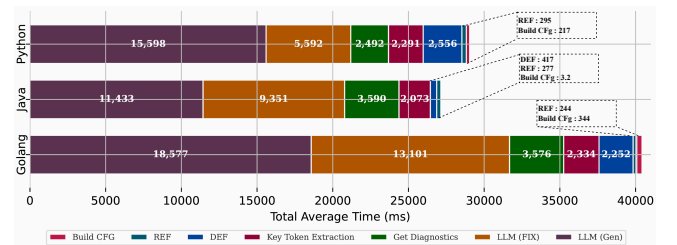
**Table 4: Overall token usage distribution and comparison.**

| Method | Java | Golang | Python | Averaged |
|---|---|---|---|---|
| *LSPRAG - Breakdown* | | | | |
| GEN | 2,144 | 3,379 | 2,674 | 2,732 |
| FIX | 3,252 | 1,458 | 586 | 1,832 |
| *Total Token Usage Comparison* | | | | |
| **LSPRAG** | 5,396 | 4,837 | 3,260 | 4,497 |
| CODEQA | 11,720 | 4,831 | 5,466 | 7,339 |
| RAG | 3,790 | 4,006 | 4,582 | 4,126 |
| DRACO | – | – | 2,439 | – |

Table 2 using the GPT4o API, following the same experimental setup as in §7.2. For each focal method, we measured the time and tokens consumed across all processes of the LSPRAG pipeline, including a maximum of five self-correction attempts. The results of latency, averaged by programming language, are detailed in Figure 6.

On average, our tool requires 28.27 seconds and 4,497 tokens per focal method to generate and refine a unit test. As detailed in 6, the majority of this time (approximately 70%) is spent on LLM API querying for generation and refinement. Our retrieval strategy, which includes key token extraction and leveraging the LSP for reference and definition providers, accounts for about 5 seconds of the total time. While the refinement stage constitutes over 30% of the generation time, we consider this a necessary and acceptable trade-off for the improvements it yields in both line coverage and the rate of valid tests.

As shown in Table 4, LSPRAG averagely uses 2,732 tokens for the initial generation and 1,832 for fixing. While LSPRAG's token consumption is not the lowest in all cases, we consider the gains in coverage justify the trade-off. For instance, compared to CODEQA, LSPRAG uses 39% fewer tokens while achieving a 135% improvement in coverage. When measured against DRACO, LSPRAG uses 9.6% more tokens for generation and 24% more for fixing, but these increases yield coverage improvements of 30% and an additional 6%, respectively. Furthermore, LSPRAG demonstrates its efficiency against a RAG approach, delivering a 174% coverage improvement with an 8.9% increase in total tokens used for generation and fixing.



**Figure 6: Test generation cost distribution per focal method, classified by programming languages.**

## 7.4 Ablation Study

To answer RQ3, we conducted an ablation study to isolate and quantify the contribution of each of LSPRAG's core components. We evaluated three configurations, with same setup from §7.2:

**Table 5: Comparison of valid rate and line coverage across different versions of LspRag. This shows that each component of LspRag shows meaningful improvement both on coverage and valid rate.**

| Project | Model | Coverage | | | Valid Rate | | |
|---|---|---|---|---|---|---|---|
| | | Naive | LspRag⁻ | LspRag | Naive | LspRag⁻ | LspRag |
| Black | GPT4o-m | 22.72% | 39.51% | 41.91% | 57.32% | 72.04% | 79.73% |
| | GPT4o | 23.56% | 42.30% | 48.10% | 56.58% | 77.99% | 85.15% |
| | DS-v3 | 39.74% | 53.36% | 56.89% | 75.62% | 87.02% | 88.03% |
| Tornado | GPT4o-m | 27.52% | 56.72% | 57.78% | 67.14% | 81.54% | 81.84% |
| | GPT4o | 29.55% | 60.03% | 60.90% | 54.36% | 82.00% | 84.53% |
| | DS-v3 | 46.04% | 63.70% | 64.62% | 89.56% | 90.36% | 91.13% |
| Commons CLI | GPT4o-m | 4.59% | 27.10% | 33.29% | 13.40% | 23.49% | 45.08% |
| | GPT4o | 12.71% | 23.16% | 34.69% | 32.61% | 28.59% | 48.12% |
| | DS-v3 | 6.46% | 28.77% | 37.72% | 17.20% | 31.32% | 60.52% |
| Commons CSV | GPT4o-m | 26.70% | 69.77% | 80.51% | 15.74% | 37.43% | 82.85% |
| | GPT4o | 39.16% | 76.05% | 78.33% | 35.69% | 54.45% | 90.9 % |
| | DS-v3 | 32.67% | 75.29% | 83.20% | 36.76% | 49.29% | 90.95% |
| Cobra | GPT4o-m | 1.36% | 9.91% | 23.03% | 1.19% | 7.13% | 23.88% |
| | GPT4o | 2.72% | 12.56% | 27.60% | 1.79% | 8.91% | 33.27% |
| | DS-v3 | 11.58% | 25.61% | 37.22% | 9.13% | 21.78% | 34.65% |
| Logrus | GPT4o-m | 2.32% | 11.55% | 23.71% | 3.33% | 18.86% | 34.02% |
| | GPT4o | 0.65% | 10.58% | 27.75% | 0.83% | 15.00% | 32.02% |
| | DS-v3 | 10.67% | 13.55% | 21.81% | 22.50% | 17.05% | 33.11% |

- **Naive**: A baseline using the same prompt template as LspRag but without LSP-guided context or code fixing.
- **LspRag⁻**: The baseline enhanced with our context retrieval on top of LSP features.
- **LspRag**: The complete LspRag system, incorporating both context retrieval and the real-time code fixing module.

The results, averaged over five runs and presented in Table 5, reveal two key findings. First, context retrieval is crucial for generating comprehensive tests. As shown in Table 5, adding context (LspRag⁻) substantially increases line coverage over the Naive baseline across all projects and LLMs. Specifically, LspRag⁻ improved line coverage by 26% to 1528% for Golang, 82% to 490% for Java, and 34% to 106% for Python. This demonstrates that the precise, semantically aware context provided by our context retrieval on top of the LSP module enables the LLM to understand class structures, method dependencies, and necessary initializations.

Second, the real-time fixing module is vital for maintaining correctness, especially when the context is complex. While rich context boosts coverage, it can be a double-edged sword. For instance, with a complex project like Commons-CLI, the large volume of context can distract the LLM, leading to syntactically incorrect code and a lower valid rate for LspRag⁻. The full LspRag system mitigates this issue. By adding the fixing module, LspRag improved line coverage by an additional 45% to 162% for Golang, 3% to 49% for Java, and 1% to 13% for Python over LspRag⁻. More importantly, it increased the valid rate by 59% to 273% for Golang, 67% to 121% for Java, and 0.5% to 10% for Python. The improvements for Python are less pronounced, which is an expected outcome due to the language's dynamic nature. Python's dynamic nature means that its LSP diagnostics, which rely on static analysis, cannot detect certain error classes that would cause compile-time failures in statically-typed

languages. Nevertheless, the fixing module proves its general utility by consistently enhancing both the validity and coverage of generated tests across all three languages, demonstrating that real-time correction is a crucial component for robust test generation.

## 8 Discussion

**Fault-Finding Capability.** While LspRag is designed to maximize test coverage, it does not explicitly optimize for the quality of test assertions. Since effective assertions are critical for detecting faults, we conducted an experiment to evaluate this capability. We injected faults into the source code using standard mutation tools (Pit [8] and MutPy [16]) and measured whether the generated tests could detect them. As summarized in Table 6, LspRag consistently outperforms

**Table 6: Comparison of mutation scores across all baselines.**

| Project | Model | LspRag | SymPrompt | CodeQA | StandardRAG | DraCo |
|---|---|---|---|---|---|---|
| Black | DS-v3 | 25.21% | 12.75% | 11.85% | 5.03% | 7.73% |
| | GPT4o | 30.28% | 21.04% | 18.47% | 9.69% | 17.03% |
| | GPT4o-m | 31.90% | 15.01% | 8.18% | 2.31% | 8.71% |
| Tornado | DS-v3 | 35.39% | 1.46% | 2.03% | 21.16% | 2.03% |
| | GPT4o | 30.37% | 9.11% | 7.66% | 18.16% | 21.69% |
| | GPT4o-m | 24.96% | 5.15% | 13.62% | 9.28% | 15.68% |
| Commons CLI | DS-v3 | 19.67% | 3.22% | 4.21% | 12.73% | — |
| | GPT4o | 19.01% | 1.32% | 7.19% | 0.41% | — |
| | GPT4o-m | 19.92% | 0.08% | 1.40% | 0.74% | — |
| Commons CSV | DS-v3 | 58.50% | 11.75% | 39.63% | 17.63% | — |
| | GPT4o | 50.50% | 0.25% | 9.38% | 5.63% | — |
| | GPT4o-m | 50.13% | 1.13% | 7.63% | 2.75% | — |

the baselines, suggesting that its accurate context retrieval not only boosts coverage but also enhances fault detection.

**Comparison with SBST Tool.** To contextualize the performance of our LspRag, we conducted a direct comparison with EvoSuite, a widely-used SBST tool for java. For fair comparison, EvoSuite was configured to generate unit tests per method under a fixed time budget of 28 seconds per method under test. The results are summarized in Table 7. The results indicate that LspRag

**Table 7: Comparison with EvoSuite on Java projects.**

| Project | Coverage | | Valid Rate | |
|---|---|---|---|---|
| | LspRag | EvoSuite | LspRag | EvoSuite |
| Commons-CLI | 37.72% | **56.33%** | 60.52% | **100%** |
| Commons-CSV | **83.20%** | 69.70% | 90.95% | **100%** |

achieves competitive code coverage against EvoSuite. EvoSuite consistently achieves a 100% valid generation rate. This is an expected outcome, as its test generation process is rooted in a strict code mutation strategy that inherently guarantees syntactically correct and compilable tests. In contrast, LspRag is designed with a different primary goal: to support multiple languages through a lightweight static analysis framework. This design choice involves a trade-off, where we sacrifice a guaranteed valid rate for broader applicability and scalability. Nevertheless, the result demonstrates that LspRag still delivers coverage results that are comparable to EvoSuite.

**Threats to Validity.** We identify two primary external threats to the generalizability of our approach. First, the effectiveness of LSPRAG is inherently coupled with the correctness of the underlying LSP server implementation for a given language. Our approach relies on the LSP server to provide accurate static analysis results (e.g., definitions, references, and type hierarchies). Consequently, if a programming language lacks robust LSP support, or for servers with erratic or incomplete analysis capabilities, the quality of the retrieved context would be degraded, potentially impacting the quality of the generated code. Second, the behavior of LSPRAG can be influenced by environment-specific configurations. For example, the behavior of LSP server is affected by different IDEs. We observed that for a Java project, VS Code [44] might define the class path as the project's top-level directory, whereas Cursor [11] sets it to the src/ subdirectory. Such discrepancies can affect which files are analyzed, potentially affecting the LSP's static analysis results. Our prototype accommodates these variations through configuration.

**Extensibility to Other Languages.** One of the key design goals of LSPRAG is extensibility. Its architecture is modular, comprising: (1) a language-specific module for AST parsing and key feature extraction, and (2) language-agnostic modules that leverage the LSP for context retrieval and code fixing. To adapt LSPRAG for a new programming language, manual implementation is only necessary for the first module. This task involves translating the language's specific tree-sitter based AST nodes into LSPRAG's IRs. This is a well-defined task requiring a modest implementation effort. Integrating the LSP-based capabilities is straightforward. It only requires the installation of the appropriate language server (e.g., as an IDE extension). Once installed, LSPRAG can immediately utilize the server's rich functionalities, such as "Go To Definition," without any further modification to its core logic.

## 9 Related Work

**Existing Approaches on Unit Test Generation.** Automated unit test generation has been a long-standing goal in software engineering, with two primary waves of innovation: classical search-based and symbolic techniques, and more recent LLM-based approaches. Classical techniques primarily focus on maximizing a specific criterion, most often code coverage. Search-Based Software Testing (SBST) tools [15, 23, 33, 47, 49] for Java, and tools [2, 29, 39] for Python, employ evolutionary algorithms or random testing to generate test suites. While powerful, these methods are *inappropriate for real-time*. They rely on whole-program compilation, heavy instrumentation, and numerous execution iterations—steps that are too slow and fragile when the codebase is in flux, changing every few minutes during active development. Their operation assumes a stable, buildable project state, which is often not the case in the early, exploratory stages of writing new code. More recently, LLMs have been applied to test generation, often outperforming classical tools in terms of the readability and initial quality of generated tests [2, 13, 48, 56, 61]. However, most of these works have selected an iterative loop either for error detection or coverage feedback. For instance, Altmayer et al. [2] propose coverage-guided pipelines where an LLM generates tests, which are then executed to gather coverage data; this feedback is used to prompt the LLM to cover the remaining branches. While these approaches have improved

previous work in terms of latency and coverage, they are still inappropriate for real-time scenarios. Furthermore, both classical and many LLM-based approaches are tailored for specific programming languages (e.g., Java and Python), limiting their generalizability.

**Retrieval-Augmented Generation.** RAG [28, 31] mitigates LLM's limitation of having finite training data. Initially proposed for NLP, it has been widely applied to code, notably through repository-level techniques. These use code-specific features for context retrieval, such as AST-based chunking [17, 50, 52] or program-analysis-based retrieval [6, 36, 37, 65]. Repository-level RAG techniques are representative of this effort. They used code-specific features to retrieve context, such as AST-based chunking [17, 50, 52], or program-analysis based retrieval [6, 36, 37, 65]. These techniques expose a fundamental trade-off. AST-based chunking is multi-lingual, but its reliance on superficial features leads to imprecise retrieval. In contrast, program-analysis-based retrieval achieves far greater precision by understanding the code's execution and data-flow relationships. However, this precision comes at a steep price: they require building and maintaining complex, language-specific static analyzers, thereby limiting their applicability to one specific language [7, 36, 37]. We compare LSPRAG with DRACO[7] and CODEQA[50] as they are the only open-sourced, referred code-aware RAG works.

**Neuro-Symbolic Approaches in Software Engineering.** Recent neuro-symbolic approaches combine LLMs with program analysis for enhanced code intelligence. This synergy aids code auditing and bug detection, with tools like LLMDFA [55] and DeepConstr [26] using LLMs for dataflow and constraint generation while symbolic techniques validate paths and mitigate hallucinations. For unit test generation, Tratto [12] uses a neural module to propose tokens and a symbolic module to constrain the search space based on program context. However, a common limitation is that they are tailored to a single programming language. In contrast, LSPRAG addresses this gap by leveraging the LSP to perform language-agnostic neuro-symbolic unit test generation.

## 10 Conclusion

We addressed the challenge of generating high-coverage unit tests for modern, multi-language software systems in real-time. We identified a key deficiency in existing approaches: their inability to precisely retrieve the necessary context, which hinders the generation of high-coverage and valid tests. To overcome this, we introduced LSPRAG, a novel framework that leverages the LSP to obtain precise code context. LSPRAG incorporates a hybrid analysis strategy to distill essential, branch-governing symbols from the retrieved context and a compile-free self-repair mechanism to ensure the syntactic validity of the generated tests without the overhead of compilation. Our extensive evaluation on real-world projects in Java, Python, and Golang demonstrates that LSPRAG significantly improves both line coverage and the rate of valid test generation across different programming languages and underlying LLMs.

## Acknowledgements

# References

[1] Mouna Abidi, Manel Grichi, and Foutse Khomh. 2019. Behind the scenes: developers' perception of multi-language practices. In *Proceedings of the 29th Annual International Conference on Computer Science and Software Engineering* (Toronto, Ontario, Canada) (*CASCON '19*). IBM Corp., USA, 72–81.

[2] Juan Altmayer Pizzorno and Emery D Berger. 2024. CoverUp: Coverage-Guided LLM-Based Test Generation. *arXiv e-prints* (2024), arXiv–2403.

[3] Andrea Arcuri and Lionel C. Briand. 2014. A hitchhiker's guide to statistical tests for assessing randomized algorithms in software engineering. *Software Testing, Verification and Reliability* 24, 3 (2014), 219–250. https://doi.org/10.1002/stvr.1486

[4] Moritz Beller, Georgios Gousios, Annibale Panichella, and Andy Zaidman. 2015. When, how, and why developers (do not) test in their IDEs. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*. 179–190.

[5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 1877–1901.

[6] Zeng-Hao Chen, Yan Ding, Qian Liu, Kui Wang, Jian-Guang Lou, Bei Zhou, and Lijie Zhang. 2022. CodeT: Code Generation with Pre-trained Language Models and Code-specific Tests. In *International Conference on Learning Representations (ICLR)*.

[7] Wei Cheng, Yuhan Wu, and Wei Hu. 2024. Dataflow-Guided Retrieval Augmentation for Repository-Level Code Completion. In *ACL*.

[8] Henry Coles, Thomas Laurent, Christopher Henard, Mike Papadakis, and Anthony Ventresque. 2016. PIT: a practical mutation testing tool for Java (demo). In *Proceedings of the 25th International Symposium on Software Testing and Analysis (ISSTA '16)*. ACM, 449–452. https://doi.org/10.1145/2931037.2948707

[9] GitHub Community. 2023. Discussion on GitHub Copilot API Integration. https://github.com/orgs/community/discussions/112339 Accessed: 2023-07-17.

[10] Stack Overflow Community. 2023. How to integrate GitHub Copilot's API in a Python script? https://stackoverflow.com/questions/79208670/how-to-integrate-github-copilots-api-in-a-python-script Accessed: 2023-07-17.

[11] Cursor. 2025. Cursor IDE. https://www.cursor.so/ Accessed: 2025-06-19.

[12] Erisa Daka, Giovanni Squillace, Michael Cooper, Alessandra Gorla, Gordon Fraser, Pankaj Garg, and Alberto Bacchelli. 2025. Tratto: A Neuro-Symbolic Approach to Deriving Axiomatic Test Oracles. In *Proceedings of the 34th ACM SIGSOFT International Symposium on Software Testing and Analysis*. ACM.

[13] Arghavan Moradi Dakhel, Amin Nikanjam, Vahid Majdinasab, Foutse Khomh, and Michel C. Desmarais. 2024. Effective test generation using pre-trained Large Language Models and mutation testing. *Inf. Softw. Technol.* 171, C (July 2024), 17 pages. https://doi.org/10.1016/j.infsof.2024.107468

[14] Ben Darnell and the Tornado Developers. 2025. Tornado: Python Web Framework and Asynchronous Networking Library. https://github.com/tornadoweb/tornado. Apache 2.0 licensed; originally developed at FriendFeed.

[15] Amirhossein Deljouyi, Roham Koohestani, Maliheh Izadi, and Andy Zaidman. 2025. Leveraging Large Language Models for Enhancing the Understandability of Generated Unit Tests . In *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*. IEEE Computer Society, Los Alamitos, CA, USA, 392–404. https://doi.org/10.1109/ICSE55347.2025.00032

[16] Anna Derezińska and Konrad Hałas. 2015. Improving Mutation Testing Process of Python Programs. In *Software Engineering in Intelligent Systems*. Advances in Intelligent Systems and Computing, Vol. 349. Springer, 233–242. https://doi.org/10.1007/978-3-319-18473-9_23

[17] Yushen Ding, Zhaoxuan Wang, Yu Zhang, Hang Wu, Zongjie Sun, Jian-Guang Lou, and Weizhu Chen. 2023. CodeRAG: A Repo-level RAG Framework for Code Generation. *arXiv preprint arXiv:2311.14813* (2023).

[18] OpenAI et al. 2024. GPT-4 Technical Report. arXiv:2303.08774 [cs.CL] https://arxiv.org/abs/2303.08774

[19] Apache Software Foundation. 2025. Apache Commons CLI. https://github.com/apache/commons-cli Accessed: 2025-01-15.

[20] Apache Software Foundation. 2025. Apache Commons CSV. https://github.com/apache/commons-csv Accessed: 2025-01-15.

[21] PSF (Python Software Foundation). 2025. Black. https://github.com/psf/black Accessed: 2025-01-15.

[22] Steve Francia. 2013. Cobra. https://github.com/spf13/cobra. Accessed: 2025-01-15.

[23] 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.

[24] Nat Friedman. 2021. Introducing GitHub Copilot: Your AI Pair Programmer. https://github.blog/2021-06-29-introducing-github-copilot-ai-pair-programmer/ Accessed: 2024-12-20.

[25] GitHub. 2024. Writing Tests with GitHub Copilot. https://docs.github.com/en/copilot/using-github-copilot/guides-on-using-github-copilot/writing-tests-with-github-copilot. Accessed: 2025-05-13.

[26] Gwihwan Go, Chijin Zhou, Quan Zhang, Xiazijian Zou, Heyuan Shi, and Yu Jiang. 2024. 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* (Vienna, Austria) (*ISSTA 2024*). Association for Computing Machinery, New York, NY, USA, 1338–1350. https://doi.org/10.1145/3650212.3680364

[27] Go Team at Google. 2024. Go Extension for Visual Studio Code. https://marketplace.visualstudio.com/items?itemName=golang.go. Accessed: 2025-05-23.

[28] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Wen-tsung Yih, Naman Goyal, and Wen-tau Chen. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 6769–6781.

[29] 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.

[30] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems* (Vancouver, BC, Canada) (*NIPS '20*). Curran Associates Inc., Red Hook, NY, USA, Article 793, 16 pages.

[31] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Gustavo Nogueira, Heinrich Piele, Wen-tau Chen, Wen-Tsung Yih, Timo Rocktäschel, et al. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural Information Processing Systems (NeurIPS)*, Vol. 33. 9459–9474.

[32] Wen Li, Na Meng, Li Li, and Haipeng Cai. 2021. Understanding language selection in multi-language software projects on GitHub. In *2021 IEEE/ACM 43rd International Conference on Software Engineering: Companion Proceedings (ICSE-Companion)*. IEEE, 256–257.

[33] Yun Lin, Jun Sun, Gordon Fraser, Ziheng Xiu, Ting Liu, and Jin Song Dong. 2020. 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*. 440–451.

[34] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437* (2024).

[35] Fang Liu, Yang Liu, Lin Shi, Houkun Huang, Ruifeng Wang, Zhen Yang, Li Zhang, Zhongqi Li, and Yuchi Ma. 2024. Exploring and evaluating hallucinations in llm-powered code generation. *arXiv preprint arXiv:2404.00971* (2024).

[36] Wei Liu, Ailun Yu, Daoguang Zan, Bo Shen, Wei Zhang, Haiyan Zhao, Zhi Jin, and Qianxiang Wang. 2024. GraphCoder: Enhancing Repository-Level Code Completion via Coarse-to-fine Retrieval Based on Code Context Graph. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering* (Sacramento, CA, USA) (*ASE '24*). Association for Computing Machinery, New York, NY, USA, 570–581. https://doi.org/10.1145/3691620.3695054

[37] Xiangyan Liu, Bo Lan, Zhiyuan Hu, Yang Liu, Zhicheng Zhang, Fei Wang, Michael Shieh, and Wenmeng Zhou. 2024. Codexgraph: Bridging large language models and code repositories via code graph databases. *arXiv preprint arXiv:2408.03910* (2024).

[38] Shuai Lu, Nan Duan, Hojae Han, Daya Guo, Seung-won Hwang, and Alexey Svyatkovskiy. 2022. ReACC: A Retrieval-Augmented Code Completion Framework. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, 6227–6240. https://doi.org/10.18653/V1/2022.ACL-LONG.431

[39] Stephan Lukasczyk and Gordon Fraser. 2022. Pynguin: Automated unit test generation for python. In *Proceedings of the ACM/IEEE 44th International Conference on Software Engineering: Companion Proceedings*. 168–172.

[40] Max Brunsfeld and others. 2024. Tree-sitter: An incremental parsing system for programming tools. https://tree-sitter.github.io/tree-sitter/. Accessed: 2025-05-23.

[41] Philip Mayer, Michael Kirsch, and Minh Anh Le. 2017. On multi-language software development, cross-language links and accompanying tools: a survey of professional software developers. *Journal of Software Engineering Research and Development* 5 (2017), 1–33.

[42] Microsoft. 2024. Language Server Protocol. https://github.com/microsoft/language-server-protocol Accessed: 2024-12-24.

[43] Microsoft. 2024. Python Extension for Visual Studio Code. https://marketplace.visualstudio.com/items?itemName=ms-python.python. Accessed: 2025-05-23.

[44] Microsoft. 2025. Visual Studio Code. https://code.visualstudio.com/ Accessed: 2025-06-19.

[45] Neverdecel. 2024. CodeRAG: An AI-powered tool for real-time codebase querying and augmentation using OpenAI and vector search. https://github.com/Neverdecel/CodeRAG. Accessed: 2025-05-18.

[46] Oracle. 2024. Extension Pack for Java. https://marketplace.visualstudio.com/ite ms?itemName=oracle.oracle-java. Accessed: 2025-05-23.

[47] 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.

[48] 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. *Proc. ACM Softw. Eng.* 1, FSE, Article 43 (July 2024), 21 pages. https://doi.org/10.1145/3643769

[49] Arkadii Sapozhnikov, Mitchell Olsthoorn, Annibale Panichella, Vladimir Kovalenko, and Pouria Derakhshanfar. 2024. TestSpark: IntelliJ IDEA's Ultimate Test Generation Companion. In *Proceedings of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion Proceedings.* 30–34.

[50] Sankalp Saxena. 2023. code_qa: LLM-based question-answering over code. https: //github.com/sankalp1999/code_qa. Accessed: 2025-05-23.

[51] Amazon Web Services. 2024. Amazon CodeWhisperer. https://aws.amazon.com /codewhisperer/ Accessed: 2025-07-12.

[52] Zhaowei Shou, Dayou Wang, and Zhou Wang. 2023. CodeGPT: A Pre-trained Language Model for Code Generation with AST-level Semantic Information. *arXiv preprint arXiv:2308.11345* (2023).

[53] Sirupsen. 2025. Logrus. https://github.com/sirupsen/logrus Accessed: 2025-01-15.

[54] Parth Thakkar. 2022. Copilot Explorer. https://thakkarparth007.github.io/copilot-explorer/posts/copilot-internals.html Accessed: 2024-12-20.

[55] Chengpeng Wang, Wuqi Zhang, Xiangzhe Xu, Zian Su, and Xiangyu Zhang. 2024. LLMDFA: Analyzing Dataflow in Code with Large Language Models. *arXiv preprint arXiv:2402.10754* (2024).

[56] Han Wang, Han Hu, Chunyang Chen, and Burak Turhan. 2024. Chat-like Asserts Prediction with the Support of Large Language Model. *arXiv preprint arXiv:2407.21429* (2024).

[57] 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* (Oulu, Finland) *(EASE '23).* Association for Computing Machinery, New York, NY, USA, 176–179. https://doi.org/10.1145/3593434.3593443

[58] Zejun Wang, Kaibo Liu, Ge Li, and Zhi Jin. 2024. HITS: High-coverage LLM-based Unit Test Generation via Method Slicing. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering.* 1258–1268.

[59] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (Eds.), Vol. 35. Curran Associates, Inc., 24824–24837.

[60] Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2024. Agentless: Demystifying llm-based software engineering agents. *arXiv preprint arXiv:2407.01489* (2024).

[61] Zhuokui Xie, Yinghao Chen, Chen Zhi, Shuiguang Deng, and Jianwei Yin. 2023. ChatUniTest: a ChatGPT-based automated unit test generation tool. *arXiv preprint arXiv:2305.04764* (2023).

[62] Haoran Yang, Yu Nong, Shaowei Wang, and Haipeng Cai. 2024. Multi-Language Software Development: Issues, Challenges, and Solutions. *IEEE Transactions on Software Engineering* 50, 3 (2024), 512–533. https://doi.org/10.1109/TSE.2024.335 8258

[63] Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and Weizhu Chen. 2023. Repocoder: Repository-level code completion through iterative retrieval and generation. *arXiv preprint arXiv:2303.12570* (2023).

[64] Feipeng Zhang, Wenyu Jiang, Jun Shu, Feng Zheng, Hongxin Wei, et al. 2024. On the noise robustness of in-context learning for text generation. *Advances in Neural Information Processing Systems* 37 (2024), 16569–16600.

[65] Fengji Zhang, Bei Zhou, Nan Duan, Alexey Svyatkovskiy, Luke Zettlemoyer, and Jian-Guang Lou. 2024. RepoCoder: Repository-Level Code Completion Through Iterative Retrieval and Generation. In *The Twelfth International Conference on Learning Representations (ICLR).*

[66] Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin. 2024. Codeagent: Enhancing code generation with tool-integrated agent systems for real-world repo-level coding challenges. *arXiv preprint arXiv:2401.07339* (2024).

[67] Ziyao Zhang, Yanlin Wang, Chong Wang, Jiachi Chen, and Zibin Zheng. 2024. Llm hallucinations in practical code generation: Phenomena, mechanism, and mitigation. *arXiv preprint arXiv:2409.20550* (2024).

[68] Yangyang Zhao, Alexander Serebrenik, Yuming Zhou, Vladimir Filkov, and Bogdan Vasilescu. 2017. The impact of continuous integration on other software development practices: a large-scale empirical study. In *2017 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE).* IEEE, 60–71.