# Large Language Models for Unit Testing: A Systematic Literature Review

QUANJUN ZHANG, School of Computer Science and Engineering, Nanjing University of Science and Technology, China

CHUNRONG FANG, SIQI GU, and YE SHANG, State Key Laboratory for Novel Software Technology, Nanjing University, China

ZHENYU CHEN, State Key Laboratory for Novel Software Technology, Nanjing University, China

LIANG XIAO, School of Computer Science and Engineering, Nanjing University of Science and Technology, China

Unit testing is a fundamental practice in modern software engineering, with the aim of ensuring the correctness, maintainability, and reliability of individual software components. Very recently, with the advances in Large Language Models (LLMs), a rapidly growing body of research has leveraged LLMs to automate various unit testing tasks, demonstrating remarkable performance and significantly reducing manual effort. However, due to ongoing explorations in the LLM-based unit testing field, it is challenging for researchers to understand existing achievements, open challenges, and future opportunities. This paper presents the first systematic literature review on the application of LLMs in unit testing until March 2025. We analyze 105 relevant papers from the perspectives of both unit testing and LLMs. We first categorize existing unit testing tasks that benefit from LLMs, e.g., test generation and oracle generation. We then discuss several critical aspects of integrating LLMs into unit testing research, including model usage, adaptation strategies, and hybrid approaches. We further summarize key challenges that remain unresolved and outline promising directions to guide future research in this area. Overall, our paper provides a systematic overview of the research landscape to the unit testing community, helping researchers gain a comprehensive understanding of achievements and promote future research. Our artifacts are publicly available at the GitHub repository: https://github.com/iSEngLab/AwesomeLLM4UT.

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

Additional Key Words and Phrases: Large Language Model, Automated Program Repair, LLM4APR

## 1 INTRODUCTION

Unit Testing (UT) aims to validate the correctness of individual components within a software system, and plays a crucial role in the software development and testing lifecycle [26, 27]. Nowadays, as modern software continues to evolve and support various critical industries, unit testing has become a standardized and even mandatory practice, forming the cornerstone of software quality and reliability [58]. However, conducting unit testing manually is often time-consuming and labor-intensive for developers and testing professionals. For example, prior work [16] has shown that developers typically spend more than 15% of their time writing unit tests. To alleviate this burden, researchers have devoted considerable efforts to automating the unit testing process, such as test case and assertion generation [7, 25]. Thus, unit testing has long been an active research topic, extensively studied over the past decades and continuously attracting interest from both academia and industry [16, 82, 104, 123].

Recently, **Large Language Models (LLMs)** have been increasingly applied in **Software Engineering (SE)**, fundamentally reshaping the research paradigm in the field [40, 104, 128]. Built on the Transformer architecture, these models are typically pre-trained on large-scale unlabeled corpora

to acquire generic language knowledge [101]. Benefiting from their advanced model architecture, vast parameters and extensive training datasets, LLMs have demonstrated remarkable progress in various software development and testing tasks, including code generation [52, 105, 107] and program repair [127, 129, 130]. In the domain of unit testing, the community has witnessed an explosion of studies equipped with LLMs [80, 86, 96, 109, 134], demonstrating notable advantages and indicating a promising future for further research.

However, due to the inherent complexity of unit testing and the rapid evolution of LLMs, integrating LLMs into unit testing workflows is a considerably complex undertaking, making it difficult for interested researchers to understand existing work. For example, existing LLM-based unit testing studies span a wide range of research perspectives (e.g., empirical [96, 114], technical [30, 109] studies), testing scenarios (e.g., test generation [85, 135] and oracle generation [21, 134]), and model usage paradigms (e.g., fine-tuning [86, 131] and prompting engineering [64, 65]). Despite the burgeoning interest and ongoing explorations in the field, the literature currently lacks a detailed and systematic review of the applications of LLMs in unit testing. This gap makes it challenging for researchers to understand the relationship between LLMs and their use in unit testing, and conduct follow-up research.

**This Work**. To bridge this gap, we present the first systematic literature review on the deployment of rapidly emerging LLM-based unit testing studies. Based on this, the community can gain a comprehensive understanding of existing LLM-based unit testing techniques, offering insights into their strengths, weaknesses, and research trends. We collect 105 relevant papers and conduct a comprehensive analysis from both unit testing and LLMs perspectives. From our analysis, we reveal crucial challenges and outline future research opportunities in the field. Overall, our work serves as a valuable resource for researchers and practitioners interested in navigating and advancing this rapidly developing area.

**Contributions.** To sum up, this work makes the following contributions:

- *Survey Methodology.* We conduct the first systematic literature review of 105 high-quality APR papers from 2020 to 2025 that utilize recent LLMs to tackle unit testing challenges.
- *Trend Analysis.* We perform a detailed analysis of selected studies in terms of publication trends and distribution of publication venues.
- *UT Perspective.* We conduct a comprehensive analysis from the perspective of unit testing to understand the distribution of unit testing tasks with LLMs and provide an in-depth discussion about how these tasks are solved with LLMs.
- *LLMs Perspective.* We conduct a comprehensive analysis from the perspective of LLMs to uncover the commonly-used LLMs, the types of prompt engineering, the input of the LLMs, as well as the accompanying techniques with these LLMs.
- *Challenges and Opportunities.* We highlight some crucial challenges of applying LLMs in the unit testing field and pinpoint promising directions for future research.

**Paper Organization.** Figure 1 summarizes the structure of this survey. The remainder of this paper is organized as follows. Section 2 introduces some basic concepts about unit testing and LLMs. Section 3 illustrates the survey methodology. Section 4 and Section 5 conduct the analysis from the perspectives of unit testing and LLMs, respectively. Section 6 discusses key challenges and research guidelines. Section 7 draws the conclusions.

## 2 BACKGROUND

In this section, we introduce the core concepts relevant to this work, including LLMs, unit testing, and related surveys.
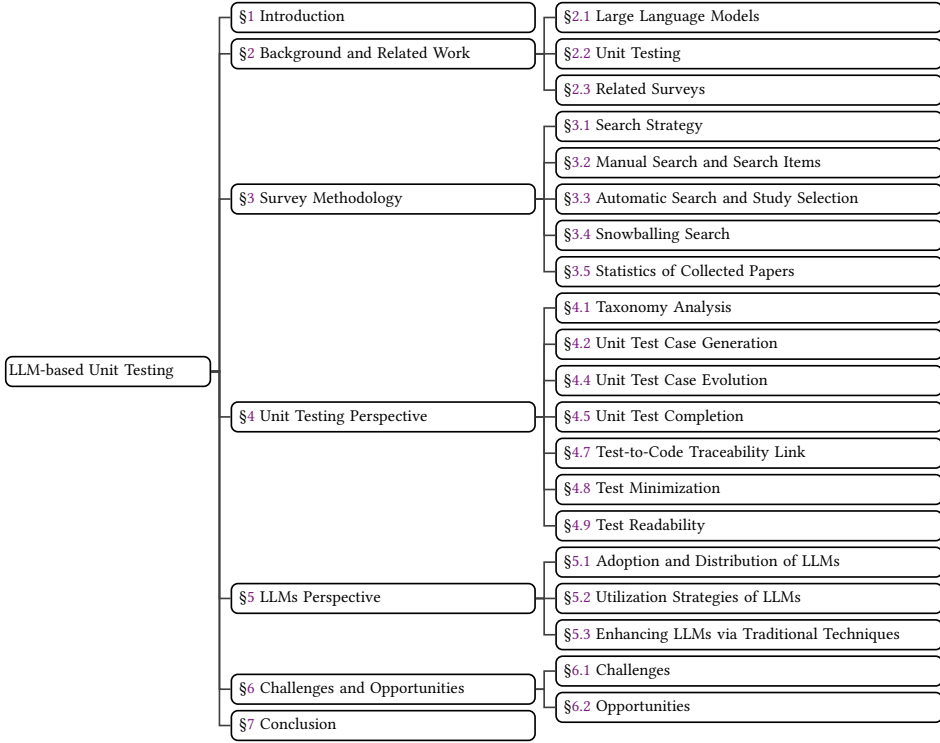
Fig. 1. Structure of this paper.

## 2.1 Large Language Models

LLMs refer to a class of large-scale Transformer-based models that are pre-trained on massive textual corpora to understand and generate human-like text [12]. To fully leverage the vast amount of unlabeled data, LLMs are typically trained using self-supervised learning objectives, such as masked language modeling [23], masked span prediction [108], and causal language modeling [67].

LLMs are primarily built on the Transformer [101] architecture, which consists of an encoder for input representation and a decoder for output generation. Based on architecture design, LLMs can be categorized into three types: (1) encoder-only models (e.g., BERT [20] and CodeBERT [23]), designed for understanding tasks; (2) encoder-decoder models (e.g., T5 [79] and CodeT5 [108]), designed for translation tasks; and (3) decoder-only models (e.g., LLaMA [98] and CodeLLaMA[81]), designed for generation tasks. In addition, based on their accessibility, LLMs can be categorized into black-box models (e.g., GPT-4 [68]), which are proprietary and closed-source, and open-source models (e.g., CodeLlama [81]), which provide public access to their architecture and weights. While commercial LLMs continue to dominate the top of the leaderboards, an increasing number of open-source models, such as CodeLlama [81] and DeepSeek-R1 [18], are emerging and demonstrating strong performance across a variety of tasks. Owing to their advanced training mechanisms, model architectures and extensive training datasets, LLMs demonstrate impressive capabilities across a wide range of SE tasks [104, 128]. In this work, we systematically investigate how LLMs have been applied to address the long-standing challenges of automating unit testing tasks.

```
//focal method:
toCollection() { return new
org.apache.commons.functor.generator.util.CollectionTransformer<E,
java.util.Collection<E>>( new java.util.ArrayList<E>());
```
```
//test prefix
testSingleStepDescending() { org.apache.commons.functor.range.FloatRange
range = org.apache.commons.functor.range.Ranges.floatRange(2.0F,
BoundType.CLOSED, (-2.0F), BoundType.OPEN, (-1.0F));
java.util.List<java.lang.Float> expected = java.util.Arrays.asList(2.0F,
1.0F, 0.0F, (-1.0F)); java.util.Collection<java.lang.Float> elements =
org.apache.commons.functor.generator.loop.IteratorToGeneratorAdapter.ada
pt(range).toCollection();
//test oracle
org.junit.Assert.assertEquals(expected, elements);
```

Fig. 2. Example of a unit test case

## 2.2 Unit Testing

Software testing plays a crucial role in SE by evaluating and ensuring the correctness, reliability, and performance of software systems [104]. Existing testing practices include unit testing, integration testing, system testing, and acceptance testing, each targeting different stages and levels in the validation process of software systems. Among them, unit testing is particularly important, as it serves as the foundation for detecting bugs at an early stage of development, thereby facilitating subsequent software development activities [58]. With its long-established history, unit testing has garnered significant attention from both academia and industry, becoming an accepted and even mandatory practice in modern software engineering.

The primary objective of unit testing is to validate individual components or units of a program by isolating them and executing unit test cases. This practice helps ensure that each component functions as expected before being integrated with other components of the system. As illustrated in Figure 2, given a focal method (i.e., the unit under test), its unit test typically consists of two components: (1) a test prefix, i.e., a sequence of statements that manipulate the unit under test to a specific state, and (2) a test oracle, i.e., an assertion that defines the expected behavior or condition to be satisfied in that state. In the literature, numerous approaches have been proposed to generate unit test cases automatically, including symbolic execution testing [11, 15], random testing [61, 73], and search-based testing [6, 25]. In addition, unit testing encompasses a variety of sub-tasks depending on the perspective, such as oracle generation [134] and test repair [116], each addressing specific challenges and scenarios. These tasks represent distinct aspects of unit testing across various levels of granularity and practical contexts, offering a comprehensive view of unit testing. Given the complexity and significance of unit testing, this work focuses on how LLMs have been leveraged to support and automate various unit testing tasks.

## 2.3 Related Surveys

There exist several surveys or literature reviews on the general application of LLMs in the broader area of software engineering [22, 40, 104, 128, 136]. Different from these studies targeting the whole software engineering/testing workflow, this work focuses specifically on the achievements of LLMs in the domain of unit testing, which remains relatively underexplored. Moreover, unit testing and its related tasks have been widely surveyed in the past [16, 82, 123]. However, these studies primarily focus on traditional unit testing techniques and were conducted prior to 2014, before the emergence of modern LLMs. Thus, they do not overlap with our research, as the first

Table 1. Comparison of related surveys and our work

| Surveys | Year | Models | Scope | STR | Time frame | # Papers |
|---------|------|--------|-------|-----|-----------|----------|
| Watson et al. [110] | 2022 | DL | Software Engineering | Y | 2009-2019 | 128 |
| Wang et al. [106] | 2022 | ML/DL | Software Engineering | Y | 2009-2020 | 1428 |
| Yang et al. [115] | 2022 | DL | Software Engineering | Y | 2015-2020 | 142 |
| Wang et al. [104] | 2023 | LLMs | Software Testing | Y | 2019-2023 | 102 |
| Fan et al. [22] | 2023 | LLMs | Software Engineering | ✗ | - | - |
| Zheng et al. [136] | 2023 | LLMs | Software Engineering | Y | 2022-2024 | 123 |
| Hou et al. [40] | 2023 | LLMs | Software Engineering | Y | 2017-2023 | 395 |
| Zhang et al. [128] | 2024 | LLMs | Software Engineering | Y | 2017-2024 | 1009 |
| Runeson et al. [82] | 2006 | - | Unit Testing | ✗ | - | - |
| Zakaria et al.[123] | 2008 | - | Unit Testing | ✓ | 2007-2009 | 27 |
| Dake et al. [16] | 2014 | - | Unit Testing | ✗ | - | - |
| Our Work | 2025 | LLMs | Unit Testing | ✓ | 2020-2025 | 105 |

LLM-based unit testing work emerged in 2020. Table 1 presents a detailed comparison between our survey and existing literature, highlighting the novelty and scope of our work. In summary, to the best of our knowledge, this is the first systematic literature review specifically focusing on the applications of LLMs for unit testing.

## 3 SURVEY METHODOLOGY

In this section, we describe our methodology for conducting this systematic literature review, including the search strategy, paper collection process, and paper trend analysis.

### 3.1 Search Strategy

Following prior SE surveys [106, 126], we adopt a three-stage "Quasi-Gold Standard" (QGS) strategy [124] to collect relevant research papers systematically. This strategy combines manual and automated search processes to construct a set of known relevant studies, which is a common practice to refine search queries and improve retrieval accuracy.

We first conduct a manual search to identify a seed set of relevant papers and derive a search string, as detailed in Section 3.2. We then utilize the search string to perform an automated search and employ a series of relatively strict filtering steps to extract the most relevant studies, as detailed in Section 3.3. We finally use a snowballing search to further complement the search results by manually inspecting references and citations, as detailed in Section 3.4. Given the large volume of relevant papers, this strategy allows us to capture the most pertinent papers while maintaining higher efficiency and rigor than a purely manual process. process. Particularly, we undertake the following three phases to search for and identify relevant studies.

### 3.2 Manual Search and Search Items

To construct the search items, we conduct a manual search from four top-tier SE conferences (ICSE, ESEC/FSE, ASE, ISSTA) and two journals (TOSEM and TSE), as listed in Table 2. The first two authors independently search these venues and propose an initial set of candidate papers involving both unit testing and LLMs. All authors then collaboratively review this collection to finalize the seed set of relevant publications. Then, we analyze the titles, abstracts, and keywords of these papers to identify search items, and conduct brainstorming to refine our search items, such as synonym substitution. Finally, this iterative process formulates the set of search items, listed as follows.

Table 2. Publication venues for manual search

| Acronym | Venues |
|---------|--------|
| ICSE | International Conference on Software Engineering |
| ESEC/FSE | Joint European Software Engineering Conference and Symposium on Foundations of Software Engineering |
| ASE | International Conference on Automated Software Engineering |
| ISSTA | International Symposium on Software Testing and Analysis |
| TOSEM | Transactions on Software Engineering Methodology |
| TSE | Transactions on Software Engineering |

- **Search items related to unit testing**: "unit testing" OR "(unit) test cases" OR "(unit) test generation" OR "(unit) test repair" OR "(unit) test evolution" OR "regression testing" OR "test refactoring" OR "test oracle" OR "test reduction" OR "test selection" OR "test readability"
- **Search items related to LLMs**: "Large Language Model(s)" OR "LLM(s)" OR "Pre-trained" OR "Pretraining" OR "PLM(s)" OR "(Code)BERT" OR "(Code)T5" OR "(Code)GPT" OR "Codex" OR "ChatGPT" OR "(Code)Llama" OR "GPT-*" OR "DeepSeek(*)" OR "Mistral"

## 3.3 Automatic Search and Study Selection

To perform the automated search, we utilize the above search items to collect relevant papers across four widely used databases, i.e., Google Scholar repository, ACM Digital Library, and IEEE Explorer Digital Library, at the end of January 2025. We restrict our search to articles published from 2017 onward, as the Transformer architecture [101], which serves as the foundation for LLMs, was introduced that year.

*3.3.1 Inclusion and Exclusion Criteria.* We define a set of inclusion and exclusion criteria to filter out papers that do not align with the scope of this survey.

**Inclusion criteria**. We define the following criteria for including papers:

- **I1**: The paper proposes a technique, tool, or framework to address unit testing tasks using LLMs.
- **I2**: The paper conducts an empirical study to evaluate LLMs in the context of unit testing.
- **I3**: The paper focuses on specific unit testing tasks (e.g., assertion generation) where LLMs are employed.

**Exclusion criteria**. We define the following criteria for excluding papers:

- **C1**: The paper does not involve any unit testing tasks, e.g., unit test generation.
- **C2**: The paper does not utilize any LLMs, e.g., only using traditional recurrent neural networks.
- **C3**: The paper only mentions LLMs in future work rather than integrating LLMs in the approach.
- **C4**: The paper is a previously published conference paper extended to a journal by the same authors.
- **C5**: The paper is not an original research study, such as literature reviews, surveys, tool demonstrations, or editorials.
- **C6**: The paper is a duplicate publication where the preprint and the published version have different titles.
- **C7**: The paper is published in a workshop or a doctoral symposium.
- **C8**: The paper is a grey publication, e.g., a technical report or thesis.
- **C9**: The paper is inaccessible in full text or not written in English.

During this process, the first two authors carefully review each paper to determine its eligibility based on the inclusion and exclusion criteria. In cases where their decisions differ, the paper is referred to the third author for a final decision. For example, following exclusion criterion **C4**, there exists one study that extends a previously published conference paper [63] to a journal version [62] by the same authors. In such cases, we retained only the extended journal version. Besides, following exclusion criterion **C6**, we identify four papers whose published versions have different titles compared to their preprints. In such cases, to avoid duplication, we retain only the published versions. After this step, we retained a total of 153 papers.

*3.3.2 Quality Assessment.* To further maximize the inclusion of high-quality papers, we design ten quality assessment questions to evaluate the relevance and rigor of included papers. For each paper, its quality is evaluated by a three-tier scoring system: criteria are rated as "yes" (1 point), "partial" (0.5 points), or "no" (0 points). If a paper accumulates a total score below the threshold of 8 points, it will be excluded from further analysis, which reduces the number of papers to 99. The designed quality assessment criteria (QAC) are listed as follows.

- **QA1**. Is the paper primarily focused on LLMs, rather than using them only as baselines?
- **QA2**. Is the paper's impact on the unit testing community explicitly stated?
- **QA3**. Are the research goals and key contributions explicitly defined?
- **QA4**. Has the paper been published in a reputable venue?
- **QA5**. Does the paper provide open-source artifacts for reproducibility, such as datasets, code, or benchmarks?
- **QA6**. Is the implementation of the proposed technique described with sufficient clarity and detail?
- **QA7** Are the experimental settings thoroughly explained, such as including hyperparameters and computing environments?
- **QA8** Are the utilized LLMs explicitly described, along with a clear explanation of how they are applied in the study?
- **QA9**. Are the evaluation metrics and results clearly aligned with research goals?
- **QA10**. Are both the contributions and limitations of the paper critically discussed?

## 3.4 Snowballing Search

To ensure the completeness of our study, we adopt a snowballing search approach [110] to manually incorporate papers that are previously overlooked yet remain pertinent to our study. Specifically, we examine all references (i.e., backward snowballing) or citations (i.e., forward snowballing) of collected papers to assess their quality and relevance to our survey. The manual inspection continues until no new relevant papers are identified, ultimately leading to the inclusion of an additional 6 papers in our survey. This rigorous procedure helps ensure that our final corpus is comprehensive and provides a solid foundation for the subsequent analysis of LLM-based unit testing techniques. Finally, we obtain 105 papers that are related to our work.

## 3.5 Statistics of Collected Papers

Figure 4 lists the number of relevant papers published across different publication venues. We find that 60% of the papers are published in peer-reviewed venues, with ICSE (11%) and TSE (8%) being the most popular conference and journal, respectively. Following them are ASE (7%), FSE (7%), and ISSTA (5%), all of which are recognized as top-tier software engineering venues. This trend indicates researchers are increasingly prioritizing the dissemination of their work through high-quality, peer-reviewed venues, which in turn drives innovation and further advances in the field. Meanwhile, around 40% of these papers, which are hosted on arXiv, have not undergone peer
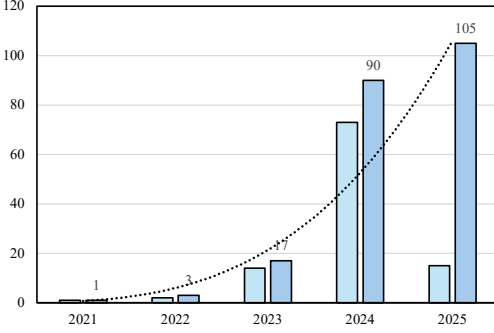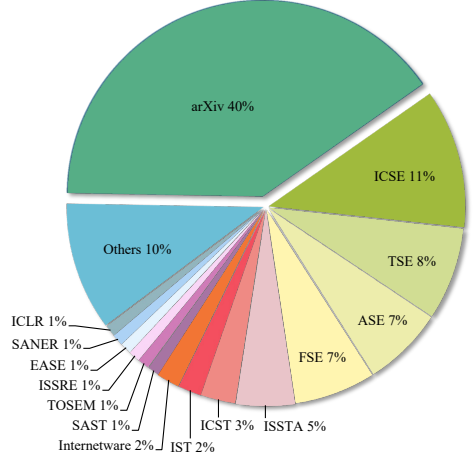
Fig. 3. Distribution of papers across years



Fig. 4. Distribution of papers across venues

review. One possible reason for this phenomenon is the sudden surge of related studies emerging within a brief period. Given the inconsistent quality of these non-peer-reviewed papers, we carry out a rigorous assessment process to ensure only high-caliber works are selected for this study.

Figure 3 further illustrates the publication trend of LLM-based unit testing. We find that the number of relevant papers appears to be rising at an almost exponential pace. While only one paper was published in 2021 and two in 2022, the number increased to 14 in 2023 and surged to 73 in 2024, with 15 already appearing by March 2025. Since our data collection ends in March 2025, the figure may not reflect the full trend for the year, and the final count is expected to increase further. These observations suggest that the application of LLMs to unit testing has gained substantial momentum since 2021 and will likely continue to be an active and growing area of research.

## 4  ANALYSIS FROM UNIT TESTING PERSPECTIVE

In this section, we present a comprehensive analysis from the perspective of unit testing and organize the collected studies according to different testing scenarios.

### 4.1  Taxonomy Analysis

Figure 5 illustrates the distribution of unit testing tasks to which LLMs are applied. It can be observed that test generation dominates the research landscape, constituting approximately 60% of the total research volume. This phenomenon is reasonable, as test case generation has long been a central component in the unit testing pipeline and continues to attract substantial attention from both academia and industry. In addition, oracle generation represents the most popular task in LLM-based unit testing research, accounting for about 14% of the research proportion. This reflects the persistent challenges in automating this task and the growing interest in leveraging LLMs to address it. Moreover, several other tasks have received moderate attention. For example, bug reproduction is explored in 5 studies (4.3%), while test evolution and test smell detection each appear in 4 studies (3.4%). Tasks such as test completion, test readability, and test minimization are addressed in 1–3 papers each, suggesting they are emerging but less mature application areas for LLMs. Notably, we observe that some studies investigate underexplored aspects of unit testing, including test-to-code traceability, test refactoring, and test validation, each appearing in only a single study. This indicates a growing research interest in broadening the scope of unit testing tasks
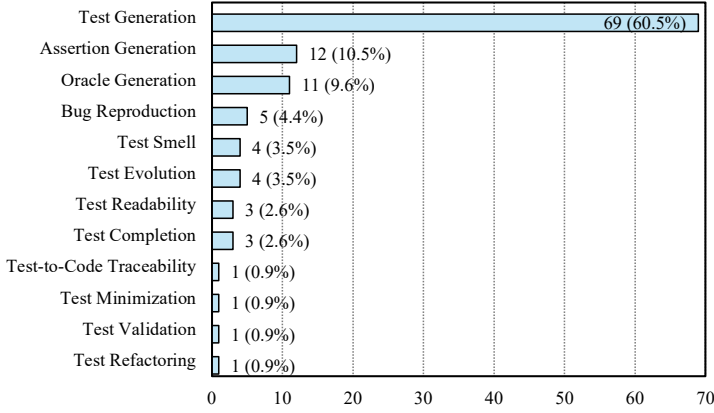
Fig. 5. Distribution of unit testing tasks with LLMs

supported by LLMs thanks to LLMs' general knowledge gleaned from vast amounts of training data. Overall, this trend highlights both the areas where LLMs have already shown strong potential and the opportunities for further exploration in more specialized or complex unit testing scenarios.

## 4.2 Unit Test Case Generation

Unit test generation typically takes a focal method (i.e., the method under test) as input and produces a complete unit test comprising two key components: a test prefix and an assertion. The test prefix sets up the required execution context, while the assertion verifies whether the focal method behaves as expected. In the literature, to reduce manual efforts in writing unit tests, researchers have proposed various automated unit test generation approaches, including symbolic execution testing [11, 15], random testing [61, 73], and search-based testing [6, 25] strategies. Traditional approaches (e.g., EvoSuite [25]) are generally designed to maximize code coverage but often struggle to produce meaningful and human-readable test cases, primarily due to their limited understanding of the semantics of focal methods. Therefore, recent research has turned to leveraging LLMs for test generation with the aim of producing more practical unit test cases from multiple dimensions, including correctness, coverage, and fault detection capability. Existing LLM-based test generation techniques can be broadly categorized into three distinct groups, discussed as follows.

*4.2.1 Training LLMs for Test Generation.* These techniques typically utilize supervised learning to train LLMs, thereby adapting them to the downstream task of unit test case generation. This is an intuitive yet effective strategy to enable LLMs to refine their pre-trained knowledge and weight parameters through limited-scale, domain-specific datasets. In the unit testing domain, the pre-training and fine-tuning paradigm has been extensively adopted during the early stages of LLM development, particularly for medium-scale models such as T5 and CodeT5, which contain hundreds of millions of parameters. The widespread adoption of this paradigm can be primarily attributed to the limited generalization capabilities of early-stage LLMs, which require task-specific training to achieve optimal performance in specialized domains like test case generation [2, 3, 99].

As early as 2020, Tufano et al. [99] introduced AthenaTest, the first LLM-based unit test generation approach, which formulates the task as a sequence-to-sequence learning problem. AthenaTest follows a two-stage training procedure: (1) denoising pre-training on a large, unsupervised Java corpus and (2) supervised fine-tuning for a downstream translation task of generating unit test

cases. To address the limitation of assertion knowledge in AthenaTest, A3Test [2] enhances BART with an assertion-aware self-supervision objective, i.e., predicting the masked tokens of a given focal method and its corresponding assert statements.

Beyond the standard pre-training and fine-tuning pipeline, researchers have explored various strategies to improve the performance of training LLMs in test case generation, including reinforcement learning [91], domain adaptation [88] and data augmentation [36]. For example, Shin et al. [88] investigate the advantages of domain adaptation for fine-tuning LLMs in the context of automated test case generation. However, most of these works rely on general-purpose, off-the-shelf LLMs that are typically trained on natural language and code corpora, without incorporating test-specific knowledge. To address this gap, Rao et al. [80] propose CAT-LM, a GPT-style LLM with 2.7 billion parameters specifically trained to learn the mapping between methods under test and their corresponding test cases, making it more suitable for the test generation task. Similarly, He et al. [34] construct a large-scale dataset, UniTSyn, containing 2.7 million method-test function pairs across five programming languages, and propose UniTester, a specialized test generation LLM trained via continual fine-tuning of SantaCoder with an autoregressive objective.

*4.2.2 Prompting LLMs for Test Case Generation.* These techniques typically construct prompts with diverse sources of project information to directly invoke LLMs for test case generation without requiring any additional training. In the unit testing community, prompt engineering has rapidly gained traction with the emergence of LLMs with billions of parameters, as their advanced capabilities make them well-suited for performing human-like interactions and generating high-quality test cases.

Most prompt-based test generation techniques primarily follow a *generation-and-refinement* paradigm, where initial test cases are first generated based on prompts and then iteratively refined using dynamic execution feedback (e.g., code coverage and failure information) to enhance their quality [13, 30, 65, 77, 83, 85, 109, 121]. For example, TestART [30] queries LLMs to generate an initial set of test cases, then uses a compiler to collect runtime information. This information is then fed into a coverage-guided testing framework and a template-based repair strategy to optimize the test cases iteratively. In addition, to design more effective prompts, researchers adopt a wide range of strategies to incorporate LLMs with valuable contextual information, including mutation testing [17], method slicing [109], demonstration retrieval [135], defect detection [119] and program analysis [113]. For example, Dakhel et al. [17] utilize mutation testing to augment the original prompt, enabling LLMs to generate tests capable of killing surviving mutants. Given a focal method, Yang et al. [113] retrieve the most similar method within the same project and incorporate it, along with its corresponding test cases, into the prompt to guide LLMs in generating test cases with project-specific exemplars.

However, most of the aforementioned work has primarily focused on standard benchmarks such as Defects4J [45] and mainstream programming languages such as Java. Recently, the research community has shifted its focus toward emerging frontier domains that present unique and diverse technical challenges. These studies include repository-level test case generation [59, 118], multi-language test case generation [76], high-performance computing software [47], industrial deployment [5, 84], Rust [14], programming problems [4], data-serialization libraries [137], game development [74], vulnerability exploitation [28], and bug reproduction [46]. Overall, these advancements demonstrate the growing versatility of prompt-based LLMs in adapting to diverse testing contexts, without the need for extensive retraining.

*4.2.3 Empirical Study.* In addition to the aforementioned technical advancements, researchers have conducted extensive empirical studies to investigate the capabilities of LLMs in generating unit test cases. These studies systematically explore the actual performance of LLMs across various

Table 3. Collected studies using LLMs for test oracle automation

| Year | Title | Oracle Type | Ref |
|------|-------|-------------|-----|
| 2020 | On Learning Meaningful Assert Statements for Unit Test Cases | Assertion Oracle | [111] |
| 2022 | TOGA: A Neural Method for Test Oracle Generation | Whole Oracle | [21] |
| 2022 | Generating Accurate Assert Statements for Unit Test Cases using Pretrained Transformers | Assertion Oracle | [100] |
| 2023 | ChatAssert: LLM-based Test Oracle Generation with External Tools Assistance | Whole Oracle | [33] |
| 2023 | Retrieval-Based Prompt Selection for Code-Related Few-Shot Learning | Assertion Oracle | [64] |
| 2023 | Using Transfer Learning for Code-Related Tasks | Assertion Oracle | [62] |
| 2023 | Neural-Based Test Oracle Generation: A Large-scale Evaluation and Lessons Learned | Whole Oracle | [38] |
| 2024 | Exploring Automated Assertion Generation via Large Language Models | Assertion Oracle | [134] |
| 2024 | AssertionBench: A Benchmark to Evaluate Large-Language Models for Assertion Generation | Assertion Oracle | [78] |
| 2024 | Chat-like Asserts Prediction with the Support of Large Language Model | Assertion Oracle | [102] |
| 2024 | Doc2Oracle: Investigating the Impact of Javadoc Comments on Test Oracle Generation | Whole Oracle | [39] |
| 2024 | Do LLMs generate test oracles that capture the actual or the expected program behaviour? | Whole Oracle | [49] |
| 2024 | Transducer Tuning: Efficient Model Adaptation for Software Tasks Using Code Property Graphs | Assertion Oracle | [122] |
| 2024 | Deep Multiple Assertions Generation | Assertion Oracle | [103] |
| 2024 | Assertify: Utilizing Large Language Models to Generate Assertions for Production Code | Assertion Oracle | [97] |
| 2025 | A Large-scale Empirical Study on Fine-tuning Large Language Models for Unit Testing | Assertion Oracle | [86] |
| 2024 | Towards More Realistic Evaluation for Neural Test Oracle Generation | Whole Oracle | [58] |
| 2024 | Assessing Evaluation Metrics for Neural Test Oracle Generation | Whole Oracle | [89] |
| 2024 | An Empirical Study on Focal Methods in Deep-Learning-Based Approaches for Assertion Generation | Assertion Oracle | [35] |
| 2025 | TOGLL: Correct and Strong Test Oracle Generation with LLMs | Whole Oracle | [37] |
| 2025 | AugmenTest: Enhancing Tests with LLM-Driven Oracles | Whole Oracle | [48] |
| 2025 | DeCon: Detecting Incorrect Assertions via Postconditions Generated by a Large Language Model | Assertion Oracle | [120] |
| 2025 | exLong: Generating Exceptional Behavior Tests with Large Language Models | Exceptional Oracle | [125] |
| 2025 | Improving Retrieval-Augmented Deep Assertion Generation via Joint Training | Assertion Oracle | [131] |
| 2025 | Improving Deep Assertion Generation via Fine-Tuning Retrieval-Augmented Pre-trained Language Models | Assertion Oracle | [132] |

aspects, including fine-tuning [86, 92], prompt engineering [53, 72, 114], integration or comparison with traditional techniques [1, 9, 44, 96, 112], source code characteristics [42], context [90], quality assessment [94], retrieval-augmented generation [87], ChatGPT-specific evaluations [31, 117, 121] and benchmarking [8, 43, 133]. For example, Shang et al. [86] conduct a large-scale empirical study to explore the potential of fine-tuning LLMs for unit testing, involving three tasks, five benchmarks, eight evaluation metrics, and 37 advanced LLMs across various architectures and sizes. Tang et al. [96] perform a systematic comparison of unit test cases generated by ChatGPT and EvoSuite based on several critical factors, including correctness, readability, code coverage, and bug detection capability. Yang et al. [114] empirically investigate the capabilities of five LLMs with various prompting settings. These empirical studies provide critical insights into the strengths and limitations of LLM-based test generation and inform future technical design and evaluation.

## 4.3 Unit Test Oracle Generation

A test oracle formally defines the expected behavior of a software unit under test for a given test prefix, serving as a critical component in unit testing. Unlike test case generation, which focuses on producing inputs and exploring execution paths, generating reliable test oracles poses a non-trivial technical challenge, as it requires capturing the intended design specification rather than merely reflecting the implemented behavior. This process demands a deep understanding of functional requirements, edge cases, and expected outcomes [39]. With the advent of LLMs, researchers have begun exploring their potential to automate oracle generation through both fine-tuning and prompting strategies. Table 3 summarizes existing LLM-based test oracle generation studies, which can be categorized into three key directions: whole oracle generation, assertion oracle generation, and exceptional oracle generation. We discuss these representative studies in detail below.

*4.3.1 Whole Oracle Generation.* In 2022, Dinella et al. [21] introduce TOGA, a transformer-based approach to infer both exceptional and assertion test oracles based on the context of the focal method. TOGA is the first LLM-powered test oracle generation study by fine-tuning CodeBERT

to (1) determine whether a test prefix raises an exception and (2) rank a set of candidate assertions. Furthermore, Liu et al. [58] identify three inappropriate settings in TOGA, i.e., generating test prefixes from correct program versions, evaluating with an unrealistic metric, and lacking a straightforward baseline. They then re-evaluate TOGA in a more realistic setting by reducing duplicates and noise during evaluation and introducing an additional ranking step to prioritize failed test cases. Similarly, Hossain et al. [38] conduct a series of replication studies to expand the understanding of TOGA's applicability, generalizability, precision, and fault detection capability across 25 real-world Java systems, 223.5K test cases, and 51K injected faults.

Unlike early-stage TOGA, which fine-tunes CodeBERT as a component for oracle classifier and ranker, recent research explores fine-tuning or prompting LLMs as core backbones to generate oracles in an end-to-end manner. For example, Hayet et al. [33] introduce CHATASSERT, a feedback-driven oracle generation technique that utilizes prompt engineering to iteratively generate and refine oracles by incorporating dynamic and static information. Khandaker et al. [48] propose AugmenTest to generate oracles for EvoSuite-generated test prefixes by prompting LLMs to infer the intended behavior of a focal method from documentation and developer comments.

In addition to the above novel techniques, researchers have conducted numerous empirical studies to explore the capabilities of LLMs in generating oracles. For example, Hossain et al. [37] present a comprehensive study by fine-tuning seven code LLMs using six distinct prompts on a large dataset consisting of 110 Java projects. Hossain et al. [39] further conduct an empirical study to investigate the impact of Javadoc comments on test oracle generation by fine-tuning ten LLMs with three different prompts. Konstantinou et al. [50] empirically investigate whether LLMs can identify the actual and expected program behavior. Besides, Shin et al [89] undertake an empirical study to reassess the performance of prior oracle generation techniques (e.g., TOGA) and ChatGPT based on both static generation metrics (e.g., BLEU and CodeBLEU) and dynamic test adequacy metrics (e.g., line coverage and mutation score). Their findings indicate a lack of statistically significant correlation between static and dynamic metrics, suggesting that existing static generation metrics do not reliably capture the quality of the generated oracles, and dynamic test adequacy metrics should serve as the principal evaluation criteria in this field.

*4.3.2 Assertion Oracle Generation.* In addition to the aforementioned whole test oracle generation work, researchers have conducted specialized explorations on assertion generation, including fine-tuning [62, 63, 100, 100, 122] and prompt engineering [64, 102]. As a seminal work in this domain, Tufano et al. [100] frame assertion generation as a sequence-to-sequence task by fine-tuning LLMs on the ATLAS dataset, where each input pair consists of a test prefix and its focal method (i.e., the method under test), and the output is an assertion. This work positions assertion generation as a downstream task for evaluating LLMs in the context of software engineering, laying the groundwork for its adoption in subsequent LLM research [62, 63, 122]. He et al. [35] conduct an empirical study on the impact of the focal method identification strategy in ATLAS and reveal its limitations in assuming the last method call before assertions as the focal method. They then introduce ATLAS+, a revised dataset where focal methods are identified using various test-to-code traceability techniques, offering a more realistic and practical evaluation framework. Unlike the above work relying on training, Nashid et al. [64] utilize few-shot learning to query Codex to generate assertion statements given focal methods and test prefixes. Wang et al. [102] introduce CLAP, which prompts LLMs to generate assertions based on chain-of-thought reasoning and iteratively refines its predictions through interactions with both LLMs and a Python interpreter.

Recently, Zhang et al. [134] conduct the first comprehensive study on fine-tuning various LLMs for automated assertion generation across two benchmarks, five LLMs, and two metrics. Their findings underscore the potential of LLM-based assertion generation to substantially alleviate the

manual workload of unit testing experts in real-world software development, thereby inspiring further research in this domain. Building on this, Zhang et al. [132] introduce RetriGen, a retrieval-augmented automated assertion generation approach, which incorporates a novel hybrid assertion retriever to refine the assertion retrieval process by leveraging both lexical and semantic similarity to identify the most relevant assertions from external codebases. Moreover, Zhang et al. [131] propose AG-RAG, which optimizes both the retriever and generator within an end-to-end pipeline using a joint training strategy, enabling them to enhance their performance through collaborative learning mutually.

*4.3.3  Exceptional Oracle Generation.* Compared to the extensive research on assertion oracle generation, there has been only one study dedicated specifically to exception oracle generation. Zhang et al. [125] introduce exLong, an LLM-based framework for automatically generating exceptional behavior test cases, i.e., checking whether the method under test throws an exception or not. Built on CodeLlama, exLong integrates program analysis to extract execution traces leading to throw statements, guard conditions, and relevant non-exceptional test cases.

## 4.4  Unit Test Case Evolution

During software evolution, source code needs to be continuously changed to satisfy new requirements or fix reported bugs. To maintain software quality, it is essential to co-evolve the corresponding unit test cases alongside the changed source code. However, this co-evolution process can pose significant challenges for developers, particularly given the constraints of limited regression testing resources and frequent releases of updated project versions. To address this, Hu et al. [41] propose CEPROT, which fine-tunes CodeT5 to update outdated test cases based on source code modifications. Given a source code change and an associated test case, CEPROT first determines whether the test case is obsolete and requires updating; if deemed obsolete, it generates a new version of the test case accordingly. Building on CEPROT, Liu et al. [56] develop SYNTER, a prompt-based approach that queries GPT-4 using contextual information to generate updated test cases. Given an obsolete test case to repair, SYNTER constructs three types of test-repair-oriented contexts via static analysis, and ranks them to select the most relevant one for guiding the repair process. In addition, Yaraghi et al. [116] introduce TARGET, a fine-tuning-based approach that adapts CodeT5+ to more realistic and executable test repair scenarios. TARGET frames test repair as a language translation task via a two-step pipeline: collecting essential contextual information that characterizes the test breakage, and utilizing the information to construct input-output training pairs for finetuning LLMs. Overall, these studies highlight the emerging role of LLMs in automating test evolution, paving the way for more robust and adaptive regression testing in evolving software systems.

## 4.5  Unit Test Completion

Test completion attempts to automatically generate the next statement within an incomplete unit test case. This task takes as input the focal method under test, the test method signature, and the preceding statements within the incomplete test body, and produces the subsequent statement to be appended. TECO [66] represents the first LLM-based test completion approach that leverages static analysis to extract code semantics and fine-tune CodeT5 to predict subsequent test statements. Furthermore, CAT-LM [80] advances the field by pre-training on both source code and test cases, followed by supervised fine-tuning, demonstrating superior performance over TECO. This line of work suggests that LLMs can effectively assist in interactive test development workflows by incrementally completing unit test cases.

## 4.6 Test Smell

Test smells refer to potential issues in unit test cases that may degrade test quality, such as brittleness, slow execution, or lack of clarity. Similar to code smells in production code, test smells suggest that a unit test may not be well-structured, robust, or efficient. Motivated by recent advances in LLMs, Lucas et al. [60] conduct an empirical study to assess the capability of LLMs in detecting test smells, demonstrating their potential to automate this process and improve testing efficiency. Furthermore, Gao et al. [29] propose UTRefactor, an LLM-based test refactoring approach to eliminate test smells and improve the quality of unit test cases. UTRefactor extracts relevant contextual information from the test code, incorporates external knowledge, and guides LLMs through a chain-of-thought process to simulate manual refactoring with improved accuracy and consistency. These studies suggest that LLMs can serve as valuable assistants in detecting and mitigating test smells, contributing to more maintainable and robust unit test cases.

## 4.7 Test-to-Code Traceability Link

Test-to-code traceability involves establishing explicit links between unit tests and the corresponding software units they are intended to validate. This practice is essential in unit testing, as it provides visibility into how tests align with production code, thereby facilitating test coverage analysis, debugging, and test maintenance. To this end, Sun et al. [93] introduce TestLinker, a hybrid approach that combines heuristic rules with LLMs to establish test-to-code traceability links at the method level. Specifically, TestLinker fine-tunes CodeT5 to learn the inherent semantic correlation between unit tests and focal methods, and then utilizes mapping rules to accurately to accurately align predicted function names with corresponding production method declarations. Despite its importance, this area remains under-explored, calling for further research on LLM-based techniques for test-to-code traceability to enable more maintainable and traceable unit test cases.

## 4.8 Test Minimization

As software evolves, unit test cases tend to grow when software evolves, making it impractical to execute all test cases with the allocated testing budgets, especially for large software systems. Test minimization attempts to improve the efficiency of unit testing by removing redundant test cases, thus reducing execution time and resource consumption while maintaining the adequacy criteria of the test suite, such as code coverage and fault detection capability. To this end, Pan et al. [75] propose LTM, a scalable and black-box test minimization approach based on LLMs and similarity analysis. LTM utilizes LLMs to extract test case embeddings and two measures to compute their similarity via two metrics, i.e., Cosine Similarity and Euclidean Distance. Based on the similarity, LTM applies a genetic algorithm to optimize the test minimization search space, which identifies the most effective subset of the original test cases within a given testing budget. However, research on LLM-based test minimization remains limited, suggesting the need for further exploration in broader regression testing scenarios such as test case prioritization and selection.

## 4.9 Test Readability

Traditional automated unit test generation techniques, particularly search-based tools (e.g., Evo-Suite), are capable of producing test cases with high code coverage. While these tools alleviate the burden of writing unit tests manually, the generated test cases often suffer from poor readability, making them difficult for developers to understand, interpret, or maintain. Thus, improving the readability of automatically generated test cases has therefore emerged as an important research direction. To address this, Delijouyi et al. [19] introduce UTGen, which integrates LLMs into the search-based test generation process, thus combining the strengths of both paradigms to generate
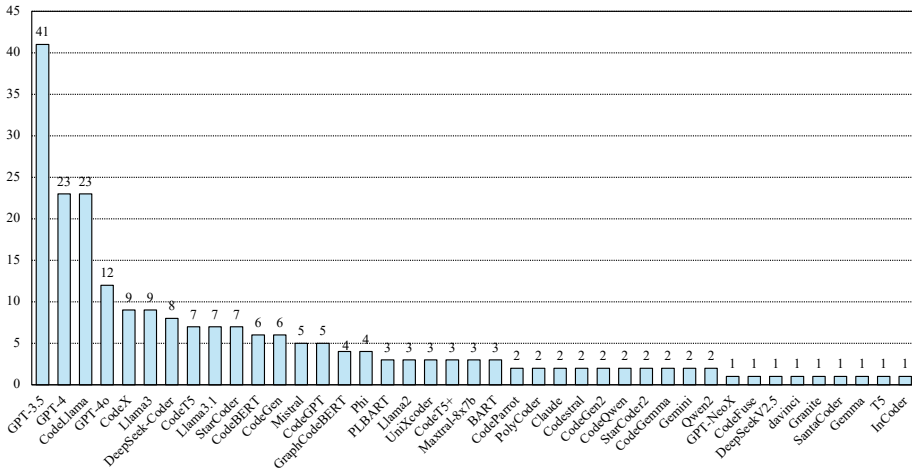
Fig. 6. Distribution of LLMs utilized in collected papers

effective and understandable test cases. UTGen utilizes LLMs to provide contextually relevant test data, insert informative comments, and suggest suitable variable names via prompt engineering. Similarly, Biagiola et al. [10] employ LLMs to improve the readability of test cases produced by EvoSuite, specifically focusing on renaming variables and test methods, while preserving functional correctness and coverage. In addition, Zhou et al. [138] introduce C3, a readability measurement tool that leverages LLMs to extract context-aware readability requirements from source code, aiming to assess and improve the readability of test inputs, especially for primitive and string types. Specifically, C3 captures the expected input context from the tested code and checks the consistency of test inputs, and its integration into EvoSuite enables the generation of more readable test inputs by guiding the test generation with these extracted contexts. These studies demonstrate the potential of LLMs in bridging the gap between traditional unit test generation techniques and human-oriented test comprehension.

> **Summary of Findings**
>
> From the perspective of unit testing, our systematic analysis reveals the following key findings: (1) LLMs have been applied across a wide range of unit testing tasks, demonstrating strong potential in both fundamental and emerging scenarios; (2) despite this breadth, research efforts remain heavily concentrated on a few primary tasks, with test case generation and oracle generation being the most studied, accounting for 20% and 19% of the collected papers, respectively; (3) although emerging tasks such as test evolution, test smell detection, and test readability are gaining attention due to LLMs' general-purpose reasoning capabilities, they remain relatively underexplored.

## 5 ANALYSIS FROM LLM PERSPECTIVE

In this section, we perform an analysis of collected papers from the viewpoint of LLMs, specifically focusing on the distribution of LLMs used, their utilization strategies, and traditional techniques employed in conjunction with LLMs.

## 5.1 Adoption and Distribution of LLMs

Figure 6 presents the distribution of LLMs adopted across the collected studies. The results reveal a clear preference for a few dominant models, particularly those accessible via commercial platforms, while open-source alternatives are gaining traction for their flexibility and adaptability.

Among all models, GPT-3.5 [69], released by OpenAI in November 2022, is the most commonly used LLM in the context of unit testing. It is trained using a combination of supervised learning and reinforcement learning from human feedback, which enables it to generate fluent, human-like responses. Primarily featured in OpenAI's ChatGPT platform, its release has significantly advanced research on LLMs in unit testing, thus contributing to its top ranking in our collected studies. Following GPT-3.5, GPT-4 [68] is the second most commonly used LLM in our collected studies. Released in March 2023, GPT-4 features a much larger parameter count and superior performance across a wide range of tasks than GPT-3.5. The third-ranked LLM is CodeLlama [81], which is an open-sourced model developed by Meta. Due to its open-source nature, researchers can not only perform prompt learning, similar to the GPT series models mentioned above, but also conduct additional training to meet domain-specific requirements. For example, Shang et al. [86] fine-tune CodeLlama to support three key unit testing tasks, i.e., test generation, test evolution, and assertion generation. In addition, GPT-4o [70], the successor to GPT-4, ranks fourth overall and has already been adopted by 12 studies since its release in May 2024. Compared with GPT-4, GPT-4o provides notable improvements in response speed, computational efficiency, and long-context reasoning, motivating some studies to explore the potential of CoT in unit testing. Other less frequently used LLMs include a range of open-source models such as DeepSeek-Coder [32], CodeT5 [108], CodeBERT [23], and StarCoder [54]. These open-source models offer flexibility for customization and experimentation, particularly in scenarios that require fine-tuning or integration with task-specific workflows.

## 5.2 Utilization Strategies of LLMs

LLMs are typically pre-trained on large-scale corpora to acquire general-purpose knowledge. Thus, a fundamental research challenge arises when integrating off-the-shelf LLMs with unit testing: *how to effectively adapt general-propose LLMs to the specialized tasks of unit testing*. In this section, we systematically categorize and analyze existing strategies for adapting LLMs to unit testing tasks.

*5.2.1 Taxonomy Analysis.* Figure 7 illustrates a hierarchical taxonomy of LLM utilization strategies in unit testing research, annotated with the number of studies adopting each approach. These adaptation strategies can be broadly grouped into two main categories: *model training* and *prompt engineering*, each comprising several specific techniques that vary in complexity, resource demands, and levels of task alignment. Overall, our analysis reveals that prompt engineering dominates the landscape with 96 studies, reflecting a growing trend toward leveraging LLMs without additional training. Zero-shot prompting (63 studies) and few-shot prompting (17 studies) are the most widely adopted strategies, while more advanced prompting techniques such as chain-of-thought (12), tree-of-thought (2), and guided tree-of-thought (2) highlight the growing sophistication of prompt engineering. On the training side, full-parameter fine-tuning remains a significant strategy, with 35 studies adapting model weights for unit testing tasks, indicating continued interest in tailoring model weights for specific unit testing objectives. Although pre-training (2 studies) and reinforcement learning (2 studies) are less commonly used, they demonstrate potential for more targeted or optimization-aware adaptation. Parameter-efficient fine-tuning (PEFT), while still emerging (4 studies), represents a practical direction toward efficient adaptation with minimal cost. In the following, we delve into each strategy in detail, discussing their technical characteristics, representative approaches, and observed trends in LLM-based unit testing research.
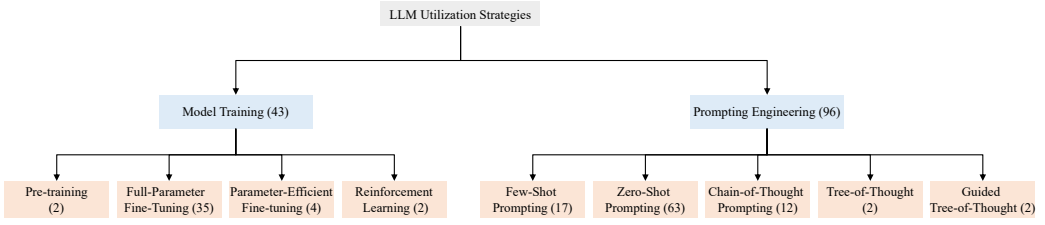
Fig. 7. Taxonomy of model utilization strategies in LLM-based unit testing research

*5.2.2   Model Training.* To support task-specific adaptation of LLMs in unit testing, various training strategies have been proposed and evaluated across the literature. These strategies differ in terms of training objectives, resource requirements, and expected performance gains. We categorize and summarize the main approaches into four types: (1) pre-training models on code-specific data, (2) full-parameter fine-tuning, (3) parameter-efficient fine-tuning (PEFT), and (4) reinforcement learning. Each of these strategies represents a distinct point in the trade-off space between model generalization, adaptation cost, and task effectiveness.

**Model Pre-training**. In the early stage, inspired by the success of LLMs in natural language, researchers attempt to adapt them to programming languages. A straightforward approach is to train code-related LLMs with similar model architectures and training strategies but using code datasets rather than natural language datasets. For example, to explore the potential of transfer learning in code-related tasks, Mastropaolo et al. [63] pre-train a T5 model with 1.5M Java methods and fine-tune it to generate assertion statements based on given focal methods and test prefixs. Furthermore, Rao et al. [80] introduce CAT-LM, a 2.7B-parameter GPT-style LLM pre-trained on aligned focal methods and test cases, making it well-suited for accurate whole-test generation.

**Full-Parameter Fine-tuning**. These techniques utilize supervised learning to train LLMs, thereby adapting them to the downstream test case generation task. This is an intuitive and effective way to allow LLMs to refine their pre-trained knowledge representations and weight parameters through limited-scale, domain-specific datasets. In the UT community, the training paradigm is extensively adopted during the early stages of LLMs with millions of parameters, such as T5 and CodeT5. The prevalent implementation of this paradigm can be primarily attributed to the limited generalization capabilities of early-stage LLMs, which require task-specific training to achieve optimal performance in specialized domains like test case generation and assertion generation [2, 3, 99]. We observe that out of 30 studies, LLMs are fine-tuned using full parameter fine-tuning techniques to adapt to downstream unit testing tasks. For example, AthenaTest [99] is the first attempt at fine-tuning LLMs for test generation, inspiring a multitude of subsequent studies [2, 88]. Similar fine-tuning-based studies on other tasks include test evolution [41, 116], assertion generation [35], and test completion [66]. Besides, Zhang et al. [134] empirically explore the potential of fine-tuning LLMs in generating assertions, and Shang et al. [86] conduct an empirical study on fine-tuning LLMs to support three unit tasks, i.e., test generation, test evolution, and assertion generation.

**Parameter-Efficient Fine-tuning (PEFT)**. However, it is quite computationally expensive and resource-intensive to fine-tune LLMs, particularly for models with billions of parameters, which usually demand substantial GPU resources. To this end, researchers introduce several parameter-efficient fine-tuning strategies in the domain of unit testing, including low-rank adaptation [8, 39, 125], (IA)$^3$ [92] and prompt tuning [92]. These studies aim to adapt LLMs with minimal weight updates, substantially reducing training costs. For example, Storhaug et al. [92] conduct the first

empirical study to extensively investigate the effectiveness of PEFT strategies, including LoRA, (IA)[3], and prompt tuning, for LLM-based unit test generation. The results demonstrate that PEFT techniques, particularly LoRA, achieve comparable performance to full-parameter fine-tuning while significantly reducing computational costs, making task-specific adaptation of LLMs more feasible. Such approaches make LLM-based unit testing more accessible and practical, particularly for academic or industrial settings with limited GPU resources.

**Reinforcement Learning (RL)**. In addition to improving fine-tuning efficiency through PEFT, some researchers also employ reinforcement learning to enhance the effectiveness of LLM-based unit testing research [91, 95]. For example, Steenhoek et al. [91] design a lightweight static analysis-based reward mode to analyze the quality of LLM-generated test cases. They then utilize the reward mode to guide the reinforcement learning process, optimizing LLMs to generate unit tests that maximize expected reward across five code quality metrics. In addition, Takerngsaksiri et al. [95] propose PyTester, an RL-based test generation approach for Python in test-driven development settings. PyTester is optimized with proximal policy optimization with a reward function that incorporates three feedback signals: syntax correctness, test executability, and code coverage.

*5.2.3 Prompt Engineering.* In addition to model training, prompt engineering has become a widely adopted strategy to guide LLM behavior for unit testing tasks. These strategies range from simple zero-shot prompts to structured, multi-agent reasoning frameworks. We categorize existing prompting strategies in the context of unit testing as follows.

**Few-shot Prompting**. Few-shot prompting provides LLMs with several input-output pairs (i.e., demonstrations) before prompting them to generate a response for new, unseen queries. For example, in assertion generation, Nashid et al. [64] retrieve similar focal methods along with their corresponding assertions based on embedding-based or frequency-based similarity, using them as in-context examples for Codex. In addition, for test generation, to provide high-quality and diverse examples, Ni et al. [65] cluster all candidate examples based on semantic similarity and select one representative from each cluster. The final prompt is then constructed with five examples, which are ordered by three different strategies: random, ascending, and descending order of cosine similarity with the target focal method.

**Zero-shot Prompting**. Zero-shot prompting directly queries LLMs to handle unseen tasks without providing any explicit examples. This prompt strategy entirely relies on LLMs' pre-existing knowledge and reasoning capabilities, thus requiring a powerful foundation model (such as Chat-GPT) for effective interaction. For example, Gu et al. [30] prompt ChatGPT to generate test cases by defining its role as a unit test case generator. The prompt also specifies the testing requirements for JDK 1.8 and JUnit 4, as well as a quality requirement to ensure branch coverage in the given focal method. In oracle generation, Konstantinou et al. [49] design a prompt that begins with a role-playing introduction and a task explanation, concluding with a sentence that specifies the desired format of the LLM's response.

**Chain-of-thought (CoT)**. CoT prompting attempts to improve LLMs' reasoning by decomposing complex problems into a sequence of intermediate reasoning steps. For example, Yang et al. [113] utilize CoT reasoning to structure the test generation process into two crucial stages: (1) providing LLMs with context information about the focal method to summarize its functionality, thus gaining a deep understanding of the method semantics; (2) guiding LLMs to iteratively generate divergent test cases that maximize branch coverage by incorporating counter-examples, thus allowing LLMs to reason step by step. Similarly, Wang et al. [109] instruct LLMs via a step-by-step CoT reasoning process, where LLMs sequentially decompose focal methods, generate test cases, and iteratively refine errors, leading to more effective and higher-coverage unit test generation. For example, during the method decomposition, LLMs first summarize the focal method's functionality, then
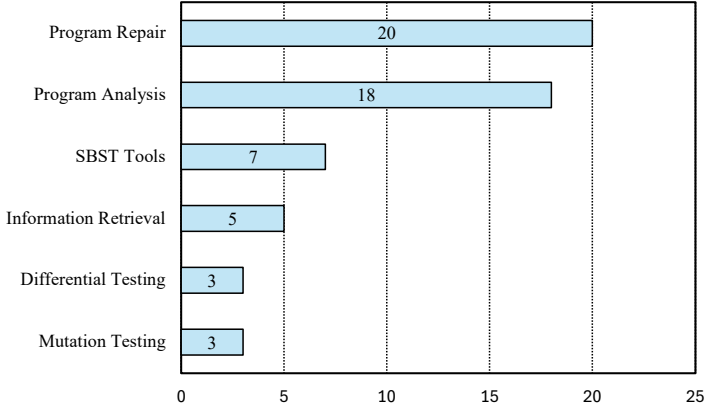
Fig. 8. Distribution of integrated traditional techniques in LLM-based unit testing research

recite the meaning and usage of all invoked non-local variables and methods, and finally break down the method into step-based meaningful slices.

**Tree-of-Thought (ToT)**. This is an advanced reasoning technique that enhances the problem-solving capabilities of LLMs by structuring reasoning as a tree-based process. Unlike CoT, which follows a linear step-by-step reasoning chain, ToT enables LLMs to iteratively explore and evaluate multiple reasoning paths, discarding suboptimal ones and converging on a more robust solution. For example, in the context of unit test generation, Ouédraogo et al. [71] employ ToT by simulating a collaborative process among three virtual software testing experts, where each expert independently proposes test cases, exchanges them with the group for evaluation, and refines them iteratively based on collaborative insights. This approach enhances test case quality by leveraging multiple reasoning pathways, ensuring coverage of typical use cases, edge cases, and exception-handling scenarios before merging the finalized test suite.

**Guided Tree-of-Thoughts (GToT)**. GToT is an enhanced version of ToT that incorporates structured, step-by-step guidance mechanisms to steer the reasoning process toward more systematic and optimal solutions. For example, Ou'edraogo et al. [71] extend ToT by explicitly instructing LLMs to follow a structured framework, including method extraction, functional test generation, edge case identification, and iterative refinement. In this setting, three virtual experts not only propose and refine test cases but also systematically analyze method signatures and exception-handling scenarios before merging their insights into a complete JUnit 4 test suite. Besides, by enforcing a well-defined reasoning structure, GToT improves test coverage, enhances code correctness, and minimizes test smells compared to the standard ToT approach.

## 5.3 Enhancing LLMs via Traditional Techniques

Although LLMs have demonstrated remarkable capabilities in various unit testing tasks, the complexity of unit testing poses unique challenges that often exceed the standalone capabilities of LLMs. To bridge this gap, recent studies have explored hybrid approaches that integrate LLMs with traditional software engineering techniques, aiming to harness the strengths of both paradigms.

*5.3.1 Taxonomy Analysis.* Figure 8 presents a taxonomy of traditional techniques integrated with LLMs in the context of unit testing, annotated with the number of studies adopting each method. These studies fall into six major categories based on the supporting techniques employed: *program*

*repair* (20 studies), *program analysis* (18 studies), *SBST tools* (7 studies), *information retrieval* (5 studies), *mutation testing* (3 studies), *differential testing* (3 studies). Among them, program repair and program analysis are the most frequently adopted techniques, indicating a strong emphasis on improving the correctness and contextual relevance of LLM-generated test cases. Program repair techniques are primarily used to fix compilation and runtime errors in generated test code, while program analysis enhances prompt relevance by extracting meaningful structural or semantic information. Other techniques, such as information retrieval and SBST tools, offer complementary support by providing dynamic context or guiding search-based exploration. Mutation and differential testing, though less frequently employed, play important roles in improving the fault-detection capabilities of generated tests. In the following, we examine each category in detail, discussing its technical motivation, representative approaches, and the observed impact on LLM-based unit testing research.

**LLMs with Program Analysis**. To facilitate LLMs' comprehension of the unit under test, an intuitive strategy is to provide them with all relevant information as prompts to the greatest extent possible. However, this strategy inevitably introduces irrelevant noise in lengthy prompts and hampers LLMs' ability to extract essential semantic information due to the lost-in-the-middle phenomenon. To this end, researchers utilize program analysis techniques to represent the unit under test and its context information more effectively, thus increasing LLMs' ability to comprehend source code accurately. For example, in the context of test generation, TELPA [113] performs backward and forward program analysis to feed LLMs with a limited number of relevant invocation methods rather than the whole focal class. Chen et al. [13] conduct program slicing to extract representative API usage and perform dependency analysis to build object construction graphs, which can provide guidance for LLM to generate meaningful test cases. Wang et al. [109] parse the focal method via static analysis to extract context information for LLMs to understand the usage and structure of the focal method. In test evolution, Yaraghi et al. [116] utilize static analysis via the Spoon library to construct method-level and class-level call graphs of the focal class, enabling the identification and prioritization of relevant repair contexts for LLMs. In assertion generation, Wang et al. [103] use Tree-sitter to convert the focal method into a data flow graph, which is then concatenated with the source code sequence as input for GraphCodeBERT.

**LLMs with Information Retrieval**. LLMs are typically trained on vast datasets, retaining learned knowledge in the form of parameters up to a fixed cutoff point. While fine-tuning provides a viable mechanism for incremental knowledge integration, frequently updating LLMs with the latest data is often impractical due to the vast number of model parameters. As a result, when generating test cases, LLMs may struggle with outdated knowledge and a lack of project-specific context. In particular, LLMs may lack awareness of new knowledge (such as updated libraries or frameworks) after their last training date and fail to incorporate critical project-level information (such as method invocation of the unit under test), which reduces the quality of generated test cases and leads to hallucinations. To address this challenge, researchers leverage information retrieval techniques to dynamically provide LLMs with relevant, up-to-date context information through prompts. For example, Nashid et al. [64] construct prompts for assertion generation by retrieving demonstrations similar to the test prefix based on embedding or frequency analysis.

**LLMs with Program Repair**. Although LLMs have achieved impressive performance on unit testing, generating syntactically correct test code remains challenging, often resulting in compilation and runtime errors. To address this, researchers have developed various strategies to identify and repair errors in generated test cases automatically. Broadly, these program repair techniques fall into two categories: dynamic feedback-based repair [66] and static pattern-based repair [2]. For example, Yuan et al. [121] adopt a validate-and-repair paradigm, iteratively refining buggy test cases by re-prompting ChatGPT with compilation error messages. In contrast, inspired by the advancements

of traditional templates program repair [130], Gu et al. [30] design five expert-informed repair templates to fix common compilation errors (such as syntax, import, and scope errors) as well as runtime errors in generated test cases.

**LLMs with Mutation Testing**. Most LLM-based unit testing approaches have demonstrated promising performance in maximizing code coverage. However, while code coverage is widely regarded as a useful metric, its correlation with actual bug detection capability remains weak. Thus, researchers employ mutation testing to simulate real bugs to optimize the bug detection ability of test cases generated by LLMs. For example, MuTAP [17] prompts LLMs to generate initial test cases via zero-shot and few-shot learning. MuTAP then applies mutation testing to assess how well the generated test cases detect faults (i.e., kill mutants). If any mutants survive, MuTAP iteratively augments the initial prompt with surviving mutants to re-prompt LLMs to generate improved test cases until all mutants are detected or no further improvements can be made. Similarly, ACH [24] attempts to construct mutants that represent faults that are both relevant to the issue at hand and undetected by existing test cases. These mutants are then used as prompts for LLMs to generate new test cases capable of detecting them.

**LLMs with Differential Testing**. Differential testing attempts to detect inconsistencies or bugs by running multiple implementations of a program with the same inputs and comparing their outputs. Given a program under test, Li et al. [55] prompt ChatGPT to infer the program's intention and ask ChatGPT to generate multiple compilable programs that share the same intention as the original. They then apply differential testing between the program under test and the ChatGPT-generated programs, using the inferred intention to identify failure-inducing test cases. Similarly, Liu et al. [57] first utilize LLMs to generate multiple program variants and test inputs. They then execute test inputs on both the program under test and its variants to construct the test oracle. Zhong et al. [137] prompt LLMs to generate high-quality unit test cases for JSON library and adopt differential testing to detect potential bugs by comparing the results from fastjson and fastjson2.

**LLMs with SBST Tools**. Another line of work combines LLMs with traditional search-based software testing (SBST) tools to leverage their respective strengths. For example, TELPA [113] utilizes the traditional search-based tool Pynguin to generate initial test cases for easily reachable branches, which are then utilized to construct prompts for LLMs, allowing them to address harder-to-cover branches further. Similarly, when Pynguin fails to increase test coverage within a given testing budget, CODAMOSA [51] invokes LLMs to generate new seed tests, which are used to resume the generation process of Pynguin. CoverUp [77] further extends CodaMosa by incorporating feedback-driven refinements based on coverage information to improve test case quality iteratively.

---

**Summary of Findings**

From the perspective of LLMs, our analysis reveals the following key findings: (1) the research landscape on LLMs for unit testing is dominated by a few commercial models, with GPT-3.5 and GPT-4 being the most widely utilized, while open-source alternatives such as CodeLlama are gaining attention for flexible adaptation; (2) prompt engineering is the most popular adaptation strategy, particularly zero-shot and few-shot prompting, although other advanced reasoning approaches, such as chain-of-thought and tree-of-thought, are emerging but remain relatively underutilized; (3) while full-parameter fine-tuning remains popular, there is increasing interest in more cost-effective alternatives, including parameter-efficient fine-tuning (e.g., LoRA) and reinforcement learning; (4) a growing number of hybrid studies that integrate LLMs with traditional techniques, such as program analysis, mutation testing, and SBST tools, highlight a trend toward improving LLMs with long-standing software engineering practices.

# 6 CHALLENGES AND OPPORTUNITIES

In this section, we discuss key challenges and highlight potential research directions across both technical and practical dimensions.

## 6.1 Challenges

Despite the promising progress in applying LLMs to unit testing, our survey reveals several persistent challenges that hinder their broader adoption and effectiveness. These challenges span from technical limitations in modeling complex software units to broader concerns around dataset quality, bug detection reliability, and model deployability. We summarize the key challenges observed in the literature as follows.

*6.1.1 Challenges in Testing Complex Units Under Test.* Unit testing aims to validate the functionality of individual software units in isolation. However, in real-world software systems, these units rarely operate independently, and they often rely on complex interactions with other functions, classes, or modules. This interconnected nature of modern codebases introduces significant challenges for LLM-based unit testing, particularly in understanding and reasoning over the broader execution context. Our analysis reveals three development stages in how contextual information is incorporated into LLMs for unit testing. From our collected papers, we observe that early studies typically feed the focal method into LLMs, prompting them to generate corresponding test cases in the form of machine translation. However, this strategy lacks critical contextual details, such as variable declarations and invoked methods, which are often necessary to produce valid and meaningful test cases. Later, with the advent of larger models featuring expanded input windows, researchers incorporate the entire focal class as additional context to capture inter-class dependencies. Although this strategy enriches the available information, it also introduces excessive input size and noise, making it difficult for LLMs to identify and focus on relevant elements. Moreover, even this expanded input fails to capture cross-file or cross-module dependencies that are common in large-scale projects. Recently, some efforts have sought to mitigate these issues by leveraging program analysis techniques to extract semantically related methods, function call relationships, and usage contexts of the unit under test. While this strategy improves the relevance of provided information, striking a balance between sufficient context and manageable input size remains an open challenge, especially in repository-level unit testing scenarios. Future research may explore more advanced static and dynamic analysis techniques (e.g., dataflow analysis) to better identify and structure relevant context. Besides, program reduction techniques, such as dead code elimination, can be employed to remove source code that is irrelevant to the behavior of the unit under test, thereby reducing complexity while preserving the essential functionality.

*6.1.2 Challenges in Detecting Real-world Software Bugs.* A fundamental purpose of unit testing is to detect individual component bugs at the early stage of software development, thereby improving software reliability. However, current evaluation metrics for unit testing primarily emphasize generation accuracy, code coverage, and defect detection capabilities, often overlooking its impact on bug detection. As illustrated in Figure 9, we summarize the distribution of evaluation metrics adopted in test generation studies based on our collected papers. It can be observed that the vast majority of studies evaluate the quality of generated test cases using code coverage (38 papers), pass rate (27 papers), and compilation rates (21 papers). There are only eight papers to identify whether the generated test cases can uncover real-world bugs.

After a careful analysis, we conclude that the challenges of existing unit testing tools in detecting real-world bugs are threefold. First, the primary reason why current studies lack sufficient focus on bug detection is the limitation of the dataset. For example, existing studies typically utilize projects
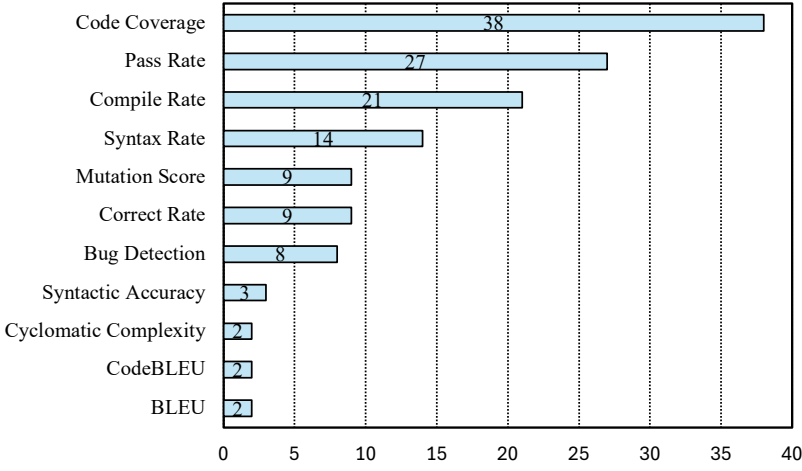
Fig. 9. Distribution of evaluation metrics in our collected test generation papers

from GitHub as evaluation datasets, which usually fail to include real-world bugs. These datasets guide researchers to evaluate the correctness of test cases (such as syntactic correctness) rather than bug detection capabilities. To address this, six papers employ mutation testing to simulate bugs for evaluation. To address this gap, we call for the development of more comprehensive and realistic datasets tailored for evaluating LLM-driven unit testing approaches. Second, oracle generation has long remained a persistent challenge in various software testing activities, including unit testing. For example, Yang et al. [114] report that GPT-4 generates valid unit test cases for only 15.74% of bugs (65/413), and among them, only 40% (39/65) are correctly detected, resulting in a final detection rate of 9.44% (39/413). Similarly, Zhang et al. [134] demonstrate that assertions generated by CodeT5 are able to detect only 15 bugs for 835 bugs from Defects4J-v2.0. The third challenge arises from the complexities inherent in real-world bug detection scenarios. Although some studies (e.g., Yang et al. [114]) attempt to evaluate bug detection, their assessment scenarios lack realism. For example, they typically generate test cases on the fixed program and run them on the buggy program. However, in practice, developers expect unit testing tools to uncover bugs in the buggy program and its fixed version is not available at that time. Besides, test cases generated by LLMs on the fixed program version may implicitly contain some information related to the patch of the bug, which should be unknown during bug finding and results in information leakage issues. Therefore, to be consistent with the real-world usage scenario, test cases should be generated on the buggy program version instead of the bug-fixed program version. A crucial challenge lies in the issue of false positives, i.e., determining whether a failing test case is caused by its assertion error or by an actual bug in the program. These false positives hinder the practical application of unit testing techniques in real-world scenarios.

*6.1.3 Challenges in Developing Unit Testing-driven LLMs.* As illustrated in Figure 6, GPT-series models are the most widely adopted in the domain of unit testing. However, the proprietary and black-box nature of these commercial models poses several challenges for real-world deployment, including concerns related to deployment environments and potential privacy breaches. For example, due to concerns regarding data privacy, when designing LLM-based unit testing tools, most organizations tend to avoid using commercial LLMs in production workflows. Instead, they often

prefer open-source alternatives that can be fine-tuned with domain-specific data in controlled environments. However, the computational cost of fine-tuning large-scale models remains prohibitive for most companies, leading them to rely on medium-sized models. However, these models often fail to match the performance of their commercial counterparts, as observed in our collected studies. Among existing LLMs, CAT-LM [80] stands out as the only model specifically designed for test case generation in relation to the method under test. However, it still relies on traditional NLP pre-training objectives and has a relatively smaller parameter scale compared to other models.

These observations point to an urgent need for the development of unit testing-oriented LLMs. Therefore, we advocate for the development of more unit test-oriented LLMs. However, developing such LLMs presents inherent challenges due to the unique characteristics of unit testing. First, while software repositories host abundant production code, the number of high-quality, labeled unit test cases remains limited, posing a significant barrier to large-scale pre-training. Second, designing pre-training objectives tailored specifically to unit testing, such as coverage maximization and assertion inference, requires further research and innovation to ensure models can generalize effectively in real-world testing scenarios.

## 6.2 Opportunities

*6.2.1 Opportunities in Developing End-to-End Testing-and-Debugging Framework.* While unit testing primarily aims to detect software bugs at the early stage of development, and program debugging focuses on automatically analyzing and fixing such bugs, the two tasks are inherently connected and mutually reinforcing. In particular, high-quality unit test cases not only expose defects, but also support key debugging activities such as root cause analysis, fault localization, patch generation, and patch validation. Despite this close relationship, prior work has largely treated unit testing and program debugging as separate research areas, limiting opportunities for integrated solutions and cross-domain enhancement.

In the future, with LLMs serving as the backbone, their powerful understanding and reasoning capabilities of LLMs present an opportunity to bridge this gap, enabling the integration of a fully automated framework for unit testing and program debugging. In such a framework, LLMs could first be used to automatically generate test cases and meaningful oracles to uncover potential defects in real-world programs. Upon identifying faults, program debugging components could be triggered to perform tasks such as root cause analysis, fault localization, and patch generation. These patches can then be validated against test cases and reintegrated into the unit testing process, forming a self-reinforcing feedback loop that continuously enhances both testing and debugging effectiveness. This direction attempts to unify two well-known research domains that are often developed in isolation into an interactive pipeline thereby benefiting the whole software development lifecycle. More importantly, such integration not only broadens the application scope of these two research areas but also boosts the capabilities of recent LLMs in advancing software quality assurance.

*6.2.2 Opportunities in Exploring More Unit Testing Tasks.* As discussed in Section 2.2, a majority of collected papers are concentrated on a limited number of unit testing scenarios, particularly test case generation and oracle generation. However, several important unit testing scenarios remain largely underexplored in the context of LLMs. For example, there is only one paper on the use of LLMs in test minimization, and still no research on test prioritization and test selection. These tasks typically differ from more widely studied ones, posing unique challenges when integrated with LLMs. For example, test case prioritization requires analyzing and ranking tens of thousands of test cases based on various features to determine an optimal execution order. Unlike the well-explored test generation task, where LLMs address it as a machine translation task by mapping a focal method to its corresponding test code, test prioritization feeds an entire test suite into LLMs, which

far exceeds the context window limitations of LLMs. To address this issue, a promising direction is to combine LLMs with traditional unit testing techniques. Specifically, it is promising to leverage LLMs' code semantic understanding capabilities by encoding test cases into vector representations, which are then combined with similarity-based test prioritization algorithms.

*6.2.3   Opportunities in Integrating Multi-modal Context Information.* From our collected papers, in the early stage of LLM-based unit testing research [2, 99], LLMs are typically provided with the focal method under test and tasked with directly generating the corresponding test case. Later, due to the dependencies among various units, studies attempt to adopt traditional techniques, such as program analysis and information retrieval, to enrich prompts with more relevant contextual information, such as invoked functions and variable definitions. However, most existing studies remain largely focused on providing accurate code-level context, often overlooking other valuable input types, such as documentation, API references, and bug reports. Given the advanced natural language understanding capabilities of LLMs, there is great potential to extract semantic intent and behavioral expectations from these textual sources, which may not be explicitly reflected in the code. Moreover, for certain types of software, such as Android applications, the multimodal capabilities of LLMs can be further explored to process diverse inputs, such as graphical user interfaces (GUIs), to support richer and more effective unit testing. We believe that the integration of source code and other multi-model contexts (such as textual and visual information) represents a promising direction for improving the completeness and accuracy of LLM-driven unit testing research.

## 7   CONCLUSION

Unit testing is a crucial and even mandatory practice in modern software engineering, facilitating software development and maintenance. Recently, Large Language Models (LLMs) have brought transformative capabilities to the unit testing domain, yielding impressive progress and indicating substantial potential in follow-up research. In this paper, we conduct the first systematic literature review of LLM-based unit testing, covering publications from 2020 to March 2025. We analyze 105 relevant studies from two complementary perspectives: the unit testing dimension, which includes tasks such as test generation, oracle generation, and test evolution; and the LLM dimension, which examines model usage, adaptation strategies, and integration with traditional techniques. We also reveal that despite promising progress, challenges remain in areas such as testing complex units, detecting real-world bugs, and developing test-oriented LLMs. To guide future research, we highlight several research opportunities, including building end-to-end testing-and-repair pipelines, expanding support for underexplored tasks, and leveraging multimodal context information. Overall, this survey will serve as a comprehensive reference for researchers and practitioners, and help advance the development of LLM-based unit testing research.

## ACKNOWLEDGMENTS

## REFERENCES

[1]  Azat Abdullin, Pouria Derakhshanfar, and Annibale Panichella. 2025. Test Wars: A Comparative Study of SBST, Symbolic Execution, and LLM-Based Approaches to Unit Test Generation. *CoRR* abs/2501.10200 (2025). https://doi.org/10.48550/ARXIV.2501.10200 arXiv:2501.10200

[2]  Saranya Alagarsamy, Chakkrit Tantithamthavorn, and Aldeida Aleti. 2024. A3Test: Assertion-Augmented Automated Test case generation. *Inf. Softw. Technol.* 176 (2024), 107565. https://doi.org/10.1016/J.INFSOF.2024.107565

[3] Saranya Alagarsamy, Chakkrit Tantithamthavorn, Chetan Arora, and Aldeida Aleti. 2024. Enhancing Large Language Models for Text-to-Testcase Generation. *CoRR* abs/2402.11910 (2024). https://doi.org/10.48550/ARXIV.2402.11910 arXiv:2402.11910

[4] Umar Alkafaween, Ibrahim Albluwi, and Paul Denny. 2025. Automating Autograding: Large Language Models as Test Suite Generators for Introductory Programming. *Journal of Computer Assisted Learning* 41, 1 (2025), e13100. https://doi.org/10.1111/jcal.13100

[5] Nadia Alshahwan, Jubin Chheda, Anastasia Finogenova, Beliz Gokkaya, Mark Harman, Inna Harper, Alexandru Marginean, Shubho Sengupta, and Eddy Wang. 2024. Automated Unit Test Improvement using Large Language Models at Meta. In *Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering, FSE 2024, Porto de Galinhas, Brazil, July 15-19, 2024*, Marcelo d'Amorim (Ed.). ACM, 185–196. https://doi.org/10.1145/3663529.3663839

[6] Luciano Baresi and Matteo Miraz. 2010. TestFul: Automatic Unit-Test Generation for Java Classes. In *Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering - Volume 2, ICSE 2010, Cape Town, South Africa, 1-8 May 2010*, Jeff Kramer, Judith Bishop, Premkumar T. Devanbu, and Sebastián Uchitel (Eds.). ACM, 281–284. https://doi.org/10.1145/1810295.1810353

[7] Earl T. Barr, Mark Harman, Phil McMinn, Muzammil Shahbaz, and Shin Yoo. 2015. The Oracle Problem in Software Testing: A Survey. *IEEE Trans. Software Eng.* 41, 5 (2015), 507–525. https://doi.org/10.1109/TSE.2014.2372785

[8] Vaishnavi Bhargava, Rajat Ghosh, and Debojyoti Dutta. 2024. CPP-UT-Bench: Can LLMs Write Complex Unit Tests in C++? *CoRR* abs/2412.02735 (2024). https://doi.org/10.48550/ARXIV.2412.02735 arXiv:2412.02735

[9] Shreya Bhatia, Tarushi Gandhi, Dhruv Kumar, and Pankaj Jalote. 2024. Unit Test Generation using Generative AI : A Comparative Performance Analysis of Autogeneration Tools. In *LLM4CODE@ICSE*. 54–61. https://doi.org/10.1145/3643795.3648396

[10] Matteo Biagiola, Gianluca Ghislotti, and Paolo Tonella. 2024. Improving the Readability of Automatically Generated Tests using Large Language Models. *CoRR* abs/2412.18843 (2024). https://doi.org/10.48550/ARXIV.2412.18843 arXiv:2412.18843

[11] Cristian Cadar, Patrice Godefroid, Sarfraz Khurshid, Corina S. Pasareanu, Koushik Sen, Nikolai Tillmann, and Willem Visser. 2011. Symbolic Execution for Software Testing in Practice: Preliminary Sssessment. In *Proceedings of the 33rd International Conference on Software Engineering, ICSE 2011, Waikiki, Honolulu , HI, USA, May 21-28, 2011*, Richard N. Taylor, Harald C. Gall, and Nenad Medvidovic (Eds.). ACM, 1066–1071. https://doi.org/10.1145/1985793.1985995

[12] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2024. A Survey on Evaluation of Large Language Models. *ACM Trans. Intell. Syst. Technol.* 15, 3 (2024), 39:1–39:45. https://doi.org/10.1145/3641289

[13] Yinghao Chen, Zehao Hu, Chen Zhi, Junxiao Han, Shuiguang Deng, and Jianwei Yin. 2024. ChatUniTest: A Framework for LLM-Based Test Generation. In *Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering, FSE 2024, Porto de Galinhas, Brazil, July 15-19, 2024*, Marcelo d'Amorim (Ed.). ACM, 572–576. https://doi.org/10.1145/3663529.3663801

[14] Xiang Cheng, Fan Sang, Yizhuo Zhai, Xiaokuan Zhang, and Taesoo Kim. 2025. RUG: Turbo LLM for Rust Unit Test Generation. In *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*. IEEE Computer Society, Los Alamitos, CA, USA, 634–634. https://doi.org/10.1109/ICSE55347.2025.00097

[15] Vitaly Chipounov, Volodymyr Kuznetsov, and George Candea. 2011. S2E: A Platform for In-Vivo Multi-Path Software Analysis. In *Proceedings of the 16th International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS 2011, Newport Beach, CA, USA, March 5-11, 2011*, Rajiv Gupta and Todd C. Mowry (Eds.). ACM, 265–278. https://doi.org/10.1145/1950365.1950396

[16] Ermira Daka and Gordon Fraser. 2014. A Survey on Unit Testing Practices and Problems. In *25th IEEE International Symposium on Software Reliability Engineering, ISSRE 2014, Naples, Italy, November 3-6, 2014*. IEEE Computer Society, 201–211. https://doi.org/10.1109/ISSRE.2014.11

[17] 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 (2024), 107468. https://doi.org/10.1016/J.INFSOF.2024.107468

[18] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua

Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, and S. S. Li. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. *CoRR* abs/2501.12948 (2025). https://doi.org/10.48550/ARXIV.2501.12948 arXiv:2501.12948

[19] Amirhossein Deljouyi, Roham Koohestani, Maliheh Izadi, and Andy Zaidman. 2024. Leveraging Large Language Models for Enhancing the Understandability of Generated Unit Tests. *CoRR* abs/2408.11710 (2024). https://doi.org/10.48550/ARXIV.2408.11710 arXiv:2408.11710

[20] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, 4171–4186. https://doi.org/10.18653/V1/N19-1423

[21] Elizabeth Dinella, Gabriel Ryan, Todd Mytkowicz, and Shuvendu K. Lahiri. 2022. TOGA: A Neural Method for Test Oracle Generation. In *44th IEEE/ACM 44th International Conference on Software Engineering, ICSE 2022, Pittsburgh, PA, USA, May 25-27, 2022*. ACM, 2130–2141. https://doi.org/10.1145/3510003.3510141

[22] Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and Jie M. Zhang. 2023. Large Language Models for Software Engineering: Survey and Open Problems. In *IEEE/ACM International Conference on Software Engineering: Future of Software Engineering, ICSE-FoSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 31–53. https://doi.org/10.1109/ICSE-FOSE59343.2023.00008

[23] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. 2020. CodeBERT: A Pre-Trained Model for Programming and Natural Languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020 (Findings of ACL, Vol. EMNLP 2020)*, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, 1536–1547. https://doi.org/10.18653/V1/2020.FINDINGS-EMNLP.139

[24] Christopher Foster, Abhishek Gulati, Mark Harman, Inna Harper, Ke Mao, Jillian Ritchey, Hervé Robert, and Shubho Sengupta. 2025. Mutation-Guided LLM-based Test Generation at Meta. *arXiv preprint arXiv:2501.12862* (2025).

[25] Gordon Fraser and Andrea Arcuri. 2011. EvoSuite: Automatic Test Suite Generation for Object-Oriented Software. In *SIGSOFT/FSE'11 19th ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE-19) and ESEC'11: 13th European Software Engineering Conference (ESEC-13), Szeged, Hungary, September 5-9, 2011*, Tibor Gyimóthy and Andreas Zeller (Eds.). ACM, 416–419. https://doi.org/10.1145/2025113.2025179

[26] Gordon Fraser and Andrea Arcuri. 2013. Whole Test Suite Generation. *IEEE Trans. Software Eng.* 39, 2 (2013), 276–291. https://doi.org/10.1109/TSE.2012.14

[27] Gordon Fraser and Andreas Zeller. 2012. Mutation-Driven Generation of Unit Tests and Oracles. *IEEE Trans. Software Eng.* 38, 2 (2012), 278–292. https://doi.org/10.1109/TSE.2011.93

[28] Yi Gao, Xing Hu, Zirui Chen, Xiaohu Yang, and Xin Xia. 2024. Unit Test Generation for Vulnerability Exploitation in Java Third-Party Libraries. *CoRR* abs/2409.16701 (2024). https://doi.org/10.48550/ARXIV.2409.16701 arXiv:2409.16701

[29] Yi Gao, Xing Hu, Xiaohu Yang, and Xin Xia. 2024. Context-Enhanced LLM-Based Framework for Automatic Test Refactoring. *CoRR* abs/2409.16739 (2024). https://doi.org/10.48550/ARXIV.2409.16739 arXiv:2409.16739

[30] Siqi Gu, Chunrong Fang, Quanjun Zhang, Fangyuan Tian, Jianyi Zhou, and Zhenyu Chen. 2024. TestART: Improving LLM-based Unit Test via Co-evolution of Automated Generation and Repair Iteration. *CoRR* abs/2408.03095 (2024). https://doi.org/10.48550/ARXIV.2408.03095 arXiv:2408.03095

[31] Vitor Guilherme and Auri Vincenzi. 2023. An Initial Investigation of ChatGPT Unit Test Generation Capability. In *8th Brazilian Symposium on Systematic and Automated Software Testing, SAST 2023, Campo Grande, MS, Brazil, September 25-29, 2023*, Awdren L. Fontão, Débora M. B. Paiva, Hudson Borges, Maria Istela Cagnin, Patrícia Gomes Fernandes, Vanessa Borges, Silvana M. Melo, Vinicius H. S. Durelli, and Edna Dias Canedo (Eds.). ACM, 15–24. https://doi.org/10.1145/3624032.3624035

[32] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. 2024. DeepSeek-Coder: When the Large Language Model Meets Programming - The Rise of Code Intelligence. *CoRR* abs/2401.14196 (2024). https://doi.org/10.48550/ARXIV.2401.14196 arXiv:2401.14196

[33] Ishrak Hayet, Adam Scott, and Marcelo d'Amorim. 2025. ChatAssert: LLM-Based Test Oracle Generation With External Tools Assistance. *IEEE Transactions on Software Engineering* 51, 1 (2025), 305–319. https://doi.org/10.1109/TSE.2024.3519159

[34] Yifeng He, Jiabo Huang, Yuyang Rong, Yiwen Guo, Ethan Wang, and Hao Chen. 2024. UniTSyn: A Large-Scale Dataset Capable of Enhancing the Prowess of Large Language Models for Program Testing. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2024, Vienna, Austria, September*

*16-20, 2024*, Maria Christakis and Michael Pradel (Eds.). ACM, 1061–1072. https://doi.org/10.1145/3650212.3680342

[35] Yibo He, Jiaming Huang, Hao Yu, and Tao Xie. 2024. An Empirical Study on Focal Methods in Deep-Learning-Based Approaches for Assertion Generation. *Proc. ACM Softw. Eng.* 1, FSE (2024), 1750–1771. https://doi.org/10.1145/3660785

[36] Yifeng He, Jicheng Wang, Yuyang Rong, and Hao Chen. 2024. Data Augmentation by Fuzzing for Neural Test Generation. *CoRR* abs/2406.08665 (2024). arXiv:2406.08665 [cs.SE] https://arxiv.org/abs/2406.08665

[37] Soneya Binta Hossain and Matthew B. Dwyer. 2024. TOGLL: Correct and Strong Test Oracle Generation with LLMs. *CoRR* abs/2405.03786 (2024). https://doi.org/10.48550/ARXIV.2405.03786 arXiv:2405.03786

[38] Soneya Binta Hossain, Antonio Filieri, Matthew B. Dwyer, Sebastian G. Elbaum, and Willem Visser. 2023. Neural-Based Test Oracle Generation: A Large-Scale Evaluation and Lessons Learned. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2023, San Francisco, CA, USA, December 3-9, 2023*, Satish Chandra, Kelly Blincoe, and Paolo Tonella (Eds.). ACM, 120–132. https://doi.org/10.1145/3611643.3616265

[39] Soneya Binta Hossain, Raygan Taylor, and Matthew B. Dwyer. 2024. Doc2Oracle: Investigating the Impact of Javadoc Comments on Test Oracle Generation. *CoRR* abs/2412.09360 (2024). https://doi.org/10.48550/ARXIV.2412.09360 arXiv:2412.09360

[40] Xinyi Hou, Yanjie Zhao, Yue Liu, Zhou Yang, Kailong Wang, Li Li, Xiapu Luo, David Lo, John Grundy, and Haoyu Wang. 2024. Large Language Models for Software Engineering: A Systematic Literature Review. *ACM Trans. Softw. Eng. Methodol.* 33, 8 (2024), 220:1–220:79. https://doi.org/10.1145/3695988

[41] Xing Hu, Zhuang Liu, Xin Xia, Zhongxin Liu, Tongtong Xu, and Xiaohu Yang. 2023. Identify and Update Test Cases When Production Code Changes: A Transformer-Based Approach. In *38th IEEE/ACM International Conference on Automated Software Engineering, ASE 2023, Luxembourg, September 11-15, 2023*. IEEE, 1111–1122. https://doi.org/10.1109/ASE56229.2023.00165

[42] Dong Huang, Jie M. Zhang, Mingzhe Du, Mark Harman, and Heming Cui. 2024. Rethinking the Influence of Source Code on Test Case Generation. *CoRR* abs/2409.09464 (2024). https://doi.org/10.48550/ARXIV.2409.09464 arXiv:2409.09464

[43] Kush Jain, Gabriel Synnaeve, and Baptiste Rozière. 2024. TestGenEval: A Real World Unit Test Generation and Test Completion Benchmark. *CoRR* abs/2410.00752 (2024). https://doi.org/10.48550/ARXIV.2410.00752 arXiv:2410.00752

[44] Zongze Jiang, Ming Wen, Jialun Cao, Xuanhua Shi, and Hai Jin. 2024. Towards Understanding the Effectiveness of Large Language Models on Directed Test Input Generation. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering, ASE 2024, Sacramento, CA, USA, October 27 - November 1, 2024*, Vladimir Filkov, Baishakhi Ray, and Minghui Zhou (Eds.). ACM, 1408–1420. https://doi.org/10.1145/3691620.3695513

[45] René Just, Darioush Jalali, and Michael D. Ernst. 2014. Defects4J: A Database of Existing Faults to Enable Controlled Testing Studies for Java Programs. In *International Symposium on Software Testing and Analysis, ISSTA '14, San Jose, CA, USA - July 21 - 26, 2014*, Corina S. Pasareanu and Darko Marinov (Eds.). ACM, 437–440. https://doi.org/10.1145/2610384.2628055

[46] Sungmin Kang, Juyeon Yoon, and Shin Yoo. 2023. Large Language Models are Few-shot Testers: Exploring LLM-based General Bug Reproduction. In *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 2312–2323. https://doi.org/10.1109/ICSE48619.2023.00194

[47] Rabimba Karanjai, Aftab Hussain, Md. Rafiqul Islam Rabin, Lei Xu, Weidong Shi, and Mohammad Amin Alipour. 2024. Harnessing the Power of LLMs: Automating Unit Test Generation for High-Performance Computing. *CoRR* abs/2407.05202 (2024). https://doi.org/10.48550/ARXIV.2407.05202 arXiv:2407.05202

[48] Shaker Mahmud Khandaker, Fitsum Kifetew, Davide Prandi, and Angelo Susi. 2025. AugmenTest: Enhancing Tests with LLM-Driven Oracles. *arXiv preprint arXiv:2501.17461* (2025).

[49] Michael Konstantinou, Renzo Degiovanni, and Mike Papadakis. 2024. Do LLMs Generate Test Oracles That Capture the Actual or the Expected Program Behaviour? *CoRR* abs/2410.21136 (2024). https://doi.org/10.48550/ARXIV.2410.21136 arXiv:2410.21136

[50] Metin Konuk, Cem Baglum, and Ugur Yayan. 2024. Evaluation of Large Language Models for Unit Test Generation. In *2024 Innovations in Intelligent Systems and Applications Conference (ASYU)*. IEEE, 1–6.

[51] 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 *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 919–931. https://doi.org/10.1109/ICSE48619.2023.00085

[52] Jia Li, Yongmin Li, Ge Li, Zhi Jin, Yiyang Hao, and Xing Hu. 2023. SkCoder: A Sketch-based Approach for Automatic Code Generation. In *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 2124–2135. https://doi.org/10.1109/ICSE48619.2023.00179

[53] Kefan Li and Yuan Yuan. 2024. Large Language Models as Test Case Generators: Performance Evaluation and Enhancement. *CoRR* abs/2404.13340 (2024). https://doi.org/10.48550/ARXIV.2404.13340 arXiv:2404.13340

[54] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy V, Jason T. Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. 2023. Starcoder: May the Source Be with You! *Trans. Mach. Learn. Res.* 2023 (2023). https://openreview.net/forum?id=KoFOg41haE

[55] Tsz On Li, Wenxi Zong, Yibo Wang, Haoye Tian, Ying Wang, Shing-Chi Cheung, and Jeff Kramer. 2023. Nuances are the Key: Unlocking ChatGPT to Find Failure-Inducing Tests with Differential Prompting. In *38th IEEE/ACM International Conference on Automated Software Engineering, ASE 2023, Luxembourg, September 11-15, 2023*. IEEE, 14–26. https://doi.org/10.1109/ASE56229.2023.00089

[56] Jun Liu, Jiwei Yan, Yuanyuan Xie, Jun Yan, and Jian Zhang. 2024. Fix the Tests: Augmenting LLMs to Repair Test Cases with Static Collector and Neural Reranker. In *35th IEEE International Symposium on Software Reliability Engineering, ISSRE 2024, Tsukuba, Japan, October 28-31, 2024*. IEEE, 367–378. https://doi.org/10.1109/ISSRE62328.2024.00043

[57] Kaibo Liu, Yiyang Liu, Zhenpeng Chen, Jie M. Zhang, Yudong Han, Yun Ma, Ge Li, and Gang Huang. 2024. LLM-Powered Test Case Generation for Detecting Tricky Bugs. *CoRR* abs/2404.10304 (2024), arXiv–2404. https://doi.org/10.48550/ARXIV.2404.10304 arXiv:2404.10304

[58] Zhongxin Liu, Kui Liu, Xin Xia, and Xiaohu Yang. 2023. Towards More Realistic Evaluation for Neural Test Oracle Generation. In *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2023, Seattle, WA, USA, July 17-21, 2023*, René Just and Gordon Fraser (Eds.). ACM, 589–600. https://doi.org/10.1145/3597926.3598080

[59] Andrea Lops, Fedelucio Narducci, Azzurra Ragone, Michelantonio Trizio, and Claudio Bartolini. 2024. A System for Automated Unit Test Generation Using Large Language Models and Assessment of Generated Test Suites. *CoRR* abs/2408.07846 (2024). https://doi.org/10.48550/ARXIV.2408.07846 arXiv:2408.07846

[60] Keila Lucas, Rohit Gheyi, Elvys Soares, Márcio Ribeiro, and Ivan Machado. 2024. Evaluating Large Language Models in Detecting Test Smells. In *Proceedings of the 38th Brazilian Symposium on Software Engineering, SBES 2024, Curitiba, Brazil, September 30 - October 4, 2024*. 672–678. https://doi.org/10.5753/SBES.2024.3642

[61] Lei Ma, Cyrille Artho, Cheng Zhang, Hiroyuki Sato, Johannes Gmeiner, and Rudolf Ramler. 2015. GRT: Program-Analysis-Guided Random Testing (T). In *30th IEEE/ACM International Conference on Automated Software Engineering, ASE 2015, Lincoln, NE, USA, November 9-13, 2015*, Myra B. Cohen, Lars Grunske, and Michael Whalen (Eds.). IEEE Computer Society, 212–223. https://doi.org/10.1109/ASE.2015.49

[62] Antonio Mastropaolo, Nathan Cooper, David Nader-Palacio, Simone Scalabrino, Denys Poshyvanyk, Rocco Oliveto, and Gabriele Bavota. 2023. Using Transfer Learning for Code-Related Tasks. *IEEE Trans. Software Eng.* 49, 4 (2023), 1580–1598. https://doi.org/10.1109/TSE.2022.3183297

[63] Antonio Mastropaolo, Simone Scalabrino, Nathan Cooper, David Nader-Palacio, Denys Poshyvanyk, Rocco Oliveto, and Gabriele Bavota. 2021. Studying the Usage of Text-To-Text Transfer Transformer to Support Code-Related Tasks. In *43rd IEEE/ACM International Conference on Software Engineering, ICSE 2021, Madrid, Spain, 22-30 May 2021*. IEEE, 336–347. https://doi.org/10.1109/ICSE43902.2021.00041

[64] Noor Nashid, Mifta Sintaha, and Ali Mesbah. 2023. Retrieval-Based Prompt Selection for Code-Related Few-Shot Learning. In *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 2450–2462. https://doi.org/10.1109/ICSE48619.2023.00205

[65] Chao Ni, Xiaoya Wang, Liushan Chen, Dehai Zhao, Zhengong Cai, Shaohua Wang, and Xiaohu Yang. 2024. CasModaTest: A Cascaded and Model-agnostic Self-directed Framework for Unit Test Generation. *CoRR* abs/2406.15743 (2024), arXiv–2406. https://doi.org/10.48550/ARXIV.2406.15743 arXiv:2406.15743

[66] Pengyu Nie, Rahul Banerjee, Junyi Jessy Li, Raymond J. Mooney, and Milos Gligoric. 2023. Learning Deep Semantics for Test Completion. In *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 2111–2123. https://doi.org/10.1109/ICSE48619.2023.00178

[67] Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2023. CodeGen: An Open Large Language Model for Code with Multi-Turn Program Synthesis. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. https://openreview.net/forum?id=iaYcJKpY2B_

[68] OpenAI. 2023. GPT-4 Technical Report. *CoRR* abs/2303.08774 (2023). https://doi.org/10.48550/ARXIV.2303.08774 arXiv:2303.08774

[69] OpenAI. 2025. GPT-3.5. URL: https://platform.openai.com/docs/models/gpt-3-5. Lasted accessed: 2025-04-01.

[70] OpenAI. 2025. GPT-4o. URL: https://platform.openai.com/docs/models/pt-4o. Lasted accessed: 2025-04-01.

[71] Wendkûuni C. Ouédraogo, Abdoul Kader Kaboré, Haoye Tian, Yewei Song, Anil Koyuncu, Jacques Klein, David Lo, and Tegawendé F. Bissyandé. 2024. Large-scale, Independent and Comprehensive study of the power of LLMs for test case generation. *CoRR* abs/2407.00225 (2024). https://doi.org/10.48550/ARXIV.2407.00225 arXiv:2407.00225

[72] Wendkûuni C. Ouédraogo, Abdoul Kader Kaboré, Haoye Tian, Yewei Song, Anil Koyuncu, Jacques Klein, David Lo, and Tegawendé F. Bissyandé. 2024. LLMs and Prompting for Unit Test Generation: A Large-Scale Evaluation. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering, ASE 2024, Sacramento, CA, USA, October 27 - November 1, 2024*, Vladimir Filkov, Baishakhi Ray, and Minghui Zhou (Eds.). ACM, 2464–2465. https://doi.org/10.1145/3691620.3695330

[73] Carlos Pacheco, Shuvendu K. Lahiri, Michael D. Ernst, and Thomas Ball. 2007. Feedback-Directed Random Test Generation. In *29th International Conference on Software Engineering (ICSE 2007), Minneapolis, MN, USA, May 20-26, 2007*. IEEE Computer Society, 75–84. https://doi.org/10.1109/ICSE.2007.37

[74] Ciprian Paduraru, Alin Stefanescu, and Augustin Jianu. 2024. Unit Test Generation using Large Language Models for Unity Game Development. In *Proceedings of the 1st ACM International Workshop on Foundations of Applied Software Engineering for Games, FaSE4Games 2024, Porto de Galinhas, Brazil, 16 July 2024*, Yann-Gaël Guéhéneuc, Fábio Petrillo, and Cristiano Politowski (Eds.). ACM. https://doi.org/10.1145/3663532.3664466

[75] Rongqi Pan, Taher Ahmed Ghaleb, and Lionel C. Briand. 2024. LTM: Scalable and Black-Box Similarity-Based Test Suite Minimization Based on Language Models. *IEEE Trans. Software Eng.* 50, 11 (2024), 3053–3070. https://doi.org/10.1109/TSE.2024.3469582

[76] Rangeet Pan, Myeongsoo Kim, Rahul Krishna, Raju Pavuluri, and Saurabh Sinha. 2024. Multi-language Unit Test Generation using LLMs. *CoRR* abs/2409.03093 (2024). https://doi.org/10.48550/ARXIV.2409.03093 arXiv:2409.03093

[77] Juan Altmayer Pizzorno and Emery D. Berger. 2024. CoverUp: Coverage-Guided LLM-Based Test Generation. *CoRR* abs/2403.16218 (2024). https://doi.org/10.48550/ARXIV.2403.16218 arXiv:2403.16218

[78] Vaishnavi Pulavarthi, Deeksha Nandal, Soham Dan, and Debjit Pal. 2024. AssertionBench: A Benchmark to Evaluate Large-Language Models for Assertion Generation. *CoRR* abs/2406.18627 (2024). https://doi.org/10.48550/ARXIV.2406.18627 arXiv:2406.18627

[79] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach. Learn. Res.* 21 (2020), 140:1–140:67. https://jmlr.org/papers/v21/20-074.html

[80] Nikitha Rao, Kush Jain, Uri Alon, Claire Le Goues, and Vincent J. Hellendoorn. 2023. CAT-LM Training Language Models on Aligned Code And Tests. In *38th IEEE/ACM International Conference on Automated Software Engineering, ASE 2023, Luxembourg, September 11-15, 2023*. IEEE, 409–420. https://doi.org/10.1109/ASE56229.2023.00193

[81] Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton-Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. Code Llama: Open Foundation Models for Code. *CoRR* abs/2308.12950 (2023). https://doi.org/10.48550/ARXIV.2308.12950 arXiv:2308.12950

[82] Per Runeson. 2006. A Survey of Unit Testing Practices. *IEEE Softw.* 23, 4 (2006), 22–29. https://doi.org/10.1109/MS.2006.91

[83] 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 (2024), 951–971. https://doi.org/10.1145/3643769

[84] 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, ICSE Companion 2024, Lisbon, Portugal, April 14-20, 2024*. ACM, 30–34. https://doi.org/10.1145/3639478.3640024

[85] Max Schäfer, Sarah Nadi, Aryaz Eghbali, and Frank Tip. 2024. An Empirical Evaluation of Using Large Language Models for Automated Unit Test Generation. *IEEE Trans. Software Eng.* 50, 1 (2024), 85–105. https://doi.org/10.1109/TSE.2023.3334955

[86] Ye Shang, Quanjun Zhang, Chunrong Fang, Siqi Gu, Jianyi Zhou, and Zhenyu Chen. 2024. A Large-scale Empirical Study on Fine-tuning Large Language Models for Unit Testing. *CoRR* abs/2412.16620 (2024). https://doi.org/10.48550/ARXIV.2412.16620 arXiv:2412.16620

[87] Jiho Shin, Reem Aleithan, Hadi Hemmati, and Song Wang. 2024. Retrieval-Augmented Test Generation: How Far Are We? *CoRR* abs/2409.12682 (2024). https://doi.org/10.48550/ARXIV.2409.12682 arXiv:2409.12682

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

*and Analysis, ISSTA 2024, Vienna, Austria, September 16-20, 2024*, Maria Christakis and Michael Pradel (Eds.). ACM, 1211–1222. https://doi.org/10.1145/3650212.3680354

[89] Jiho Shin, Hadi Hemmati, Moshi Wei, and Song Wang. 2024. Assessing Evaluation Metrics for Neural Test Oracle Generation. *IEEE Trans. Software Eng.* 50, 9 (2024), 2337–2349. https://doi.org/10.1109/TSE.2024.3433463

[90] Mohammed Latif Siddiq, Joanna Cecilia da Silva Santos, Ridwanul Hasan Tanvir, Noshin Ulfat, Fahmid Al Rifat, and Vinicius Carvalho Lopes. 2024. 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 2024, Salerno, Italy, June 18-21, 2024.* ACM, 313–322. https://doi.org/10.1145/3661167.3661216

[91] Benjamin Steenhoek, Michele Tufano, Neel Sundaresan, and Alexey Svyatkovskiy. 2023. Reinforcement Learning from Automatic Feedback for High-Quality Unit Test Generation. *CoRR* abs/2310.02368 (2023). https://doi.org/10.48550/ARXIV.2310.02368 arXiv:2310.02368

[92] André Storhaug and Jingyue Li. 2024. Parameter-Efficient Fine-Tuning of Large Language Models for Unit Test Generation: An Empirical Study. *CoRR* abs/2411.02462 (2024). https://doi.org/10.48550/ARXIV.2411.02462 arXiv:2411.02462

[93] Weifeng Sun, Zhenting Guo, Meng Yan, Zhongxin Liu, Yan Lei, and Hongyu Zhang. 2024. Method-Level Test-to-Code Traceability Link Construction by Semantic Correlation Learning. *IEEE Trans. Software Eng.* 50, 10 (2024), 2656–2676. https://doi.org/10.1109/TSE.2024.3449917

[94] Hamed Taherkhani and Hadi Hemmati. 2024. VALTEST: Automated Validation of Language Model Generated Test Cases. *CoRR* abs/2411.08254 (2024). https://doi.org/10.48550/ARXIV.2411.08254 arXiv:2411.08254

[95] Wannita Takerngsaksiri, Rujikorn Charakorn, Chakkrit Tantithamthavorn, and Yuan-Fang Li. 2024. TDD Without Tears: Towards Test Case Generation from Requirements through Deep Reinforcement Learning. *CoRR* abs/2401.07576 (2024). https://doi.org/10.48550/ARXIV.2401.07576 arXiv:2401.07576

[96] Yutian Tang, Zhijie Liu, Zhichao Zhou, and Xiapu Luo. 2024. ChatGPT vs SBST: A Comparative Assessment of Unit Test Suite Generation. *IEEE Trans. Software Eng.* 50, 6 (2024), 1340–1359. https://doi.org/10.1109/TSE.2024.3382365

[97] Mohammad Jalili Torkamani, Abhinav Sharma, Nikita Mehrotra, and Rahul Purandare. 2024. ASSERTIFY: Utilizing Large Language Models to Generate Assertions for Production Code. *CoRR* abs/2411.16927 (2024). https://doi.org/10.48550/ARXIV.2411.16927 arXiv:2411.16927

[98] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. *CoRR* abs/2302.13971 (2023). https://doi.org/10.48550/ARXIV.2302.13971 arXiv:2302.13971

[99] Michele Tufano, Dawn Drain, Alexey Svyatkovskiy, Shao Kun Deng, and Neel Sundaresan. 2020. Unit Test Case Generation with Transformers and Focal Context. *CoRR* abs/2009.05617 (2020), arXiv–2009.

[100] Michele Tufano, Dawn Drain, Alexey Svyatkovskiy, and Neel Sundaresan. 2022. Generating Accurate Assert Statements for Unit Test Cases using Pretrained Transformers. In *IEEE/ACM International Conference on Automation of Software Test, AST@ICSE 2022, Pittsburgh, PA, USA, May 21-22, 2022.* ACM/IEEE, 54–64. https://doi.org/10.1145/3524481.3527220

[101] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 5998–6008. https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html

[102] Han Wang, Han Hu, Chunyang Chen, and Burak Turhan. 2024. Chat-like Asserts Prediction with the Support of Large Language Model. *CoRR* abs/2407.21429 (2024). https://doi.org/10.48550/ARXIV.2407.21429 arXiv:2407.21429

[103] Hailong Wang, Tongtong Xu, and Bei Wang. 2024. Deep Multiple Assertions Generation. In *Proceedings of the 2024 IEEE/ACM First International Conference on AI Foundation Models and Software Engineering, FORGE 2024, Lisbon, Portugal, 14 April 2024*, David Lo, Xin Xia, Massimiliano Di Penta, and Xing Hu (Eds.). ACM, 1–11. https://doi.org/10.1145/3650105.3652293

[104] Junjie Wang, Yuchao Huang, Chunyang Chen, Zhe Liu, Song Wang, and Qing Wang. 2024. Software Testing With Large Language Models: Survey, Landscape, and Vision. *IEEE Trans. Software Eng.* 50, 4 (2024), 911–936. https://doi.org/10.1109/TSE.2024.3368208

[105] Shangwen Wang, Mingyang Geng, Bo Lin, Zhensu Sun, Ming Wen, Yepang Liu, Li Li, Tegawendé F. Bissyandé, and Xiaoguang Mao. 2023. Natural Language to Code: How Far Are We?. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2023, San Francisco, CA, USA, December 3-9, 2023*, Satish Chandra, Kelly Blincoe, and Paolo Tonella (Eds.). ACM, 375–387. https://doi.org/10.1145/3611643.3616323

[106] Simin Wang, Liguo Huang, Amiao Gao, Jidong Ge, Tengfei Zhang, Haitao Feng, Ishna Satyarth, Ming Li, He Zhang, and Vincent Ng. 2023. Machine/Deep Learning for Software Engineering: A Systematic Literature Review. *IEEE*

*Trans. Software Eng.* 49, 3 (2023), 1188–1231. https://doi.org/10.1109/TSE.2022.3173346

[107] Shangwen Wang, Bo Lin, Zhensu Sun, Ming Wen, Yepang Liu, Yan Lei, and Xiaoguang Mao. 2023. Two Birds with One Stone: Boosting Code Generation and Code Search via a Generative Adversarial Network. *Proc. ACM Program. Lang.* 7, OOPSLA2 (2023), 486–515. https://doi.org/10.1145/3622815

[108] Yue Wang, Weishi Wang, Shafiq R. Joty, and Steven C. H. Hoi. 2021. CodeT5: Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Understanding and Generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, 8696–8708. https://doi.org/10.18653/V1/2021.EMNLP-MAIN.685

[109] 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, ASE 2024, Sacramento, CA, USA, October 27 - November 1, 2024*, Vladimir Filkov, Baishakhi Ray, and Minghui Zhou (Eds.). ACM, 1258–1268. https://doi.org/10.1145/3691620.3695501

[110] Cody Watson, Nathan Cooper, David Nader-Palacio, Kevin Moran, and Denys Poshyvanyk. 2022. A Systematic Literature Review on the Use of Deep Learning in Software Engineering Research. *ACM Trans. Softw. Eng. Methodol.* 31, 2 (2022), 32:1–32:58. https://doi.org/10.1145/3485275

[111] Cody Watson, Michele Tufano, Kevin Moran, Gabriele Bavota, and Denys Poshyvanyk. 2020. On Learning Meaningful Assert Statements for Unit Test Cases. In *ICSE '20: 42nd International Conference on Software Engineering, Seoul, South Korea, 27 June - 19 July, 2020*, Gregg Rothermel and Doo-Hwan Bae (Eds.). ACM, 1398–1409. https://doi.org/10.1145/3377811.3380429

[112] Danni Xiao, Yimeng Guo, Yanhui Li, and Lin Chen. 2024. Optimizing Search-Based Unit Test Generation with Large Language Models: An Empirical Study. In *Proceedings of the 15th Asia-Pacific Symposium on Internetware, Internetware 2024, Macau, SAR, China, July 24-26, 2024*, Hong Mei, Jian Lv, Abdelsalam Helal, Xiaoxing Ma, Shing-Chi Cheung, Jie Zhang, and Tao Zhang (Eds.). ACM. https://doi.org/10.1145/3671016.3674813

[113] Chen Yang, Junjie Chen, Bin Lin, Jianyi Zhou, and Ziqi Wang. 2024. Enhancing LLM-based Test Generation for Hard-to-Cover Branches via Program Analysis. *CoRR* abs/2404.04966 (2024), arXiv–2404. https://doi.org/10.48550/ARXIV.2404.04966 arXiv:2404.04966

[114] Lin Yang, Chen Yang, Shutao Gao, Weijing Wang, Bo Wang, Qihao Zhu, Xiao Chu, Jianyi Zhou, Guangtai Liang, Qianxiang Wang, and Junjie Chen. 2024. On the Evaluation of Large Language Models in Unit Test Generation. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering, ASE 2024, Sacramento, CA, USA, October 27 - November 1, 2024*, Vladimir Filkov, Baishakhi Ray, and Minghui Zhou (Eds.). ACM, 1607–1619. https://doi.org/10.1145/3691620.3695529

[115] Yanming Yang, Xin Xia, David Lo, and John C. Grundy. 2022. A Survey on Deep Learning for Software Engineering. *ACM Comput. Surv.* 54, 10s (2022), 206:1–206:73. https://doi.org/10.1145/3505243

[116] Ahmadreza Saboor Yaraghi, Darren Holden, Nafiseh Kahani, and Lionel C. Briand. 2024. Automated Test Case Repair Using Language Models. *CoRR* abs/2401.06765 (2024). https://doi.org/10.48550/ARXIV.2401.06765 arXiv:2401.06765

[117] Gaolei Yi, Zizhao Chen, Zhenyu Chen, W. Eric Wong, and Nicholas Chau. 2023. Exploring the Capability of ChatGPT in Test Generation. In *23rd IEEE International Conference on Software Quality, Reliability, and Security, QRS 2023 Companion, Chiang Mai, Thailand, October 22-26, 2023*. IEEE, 72–80. https://doi.org/10.1109/QRS-C60940.2023.00013

[118] Xin Yin, Chao Ni, Xinrui Li, Liushan Chen, Guojun Ma, and Xiaohu Yang. 2025. Enhancing LLM's Ability to Generate More Repository-Aware Unit Tests Through Precise Contextual Information Injection. *CoRR* abs/2501.07425 (2025).

[119] Xin Yin, Chao Ni, Xiaodan Xu, and Xiaohu Yang. 2024. What You See Is What You Get: Attention-based Self-guided Automatic Unit Test Generation. *CoRR* abs/2412.00828 (2024). https://doi.org/10.48550/ARXIV.2412.00828 arXiv:2412.00828

[120] Hao Yu, Tianyu Chen, Jiaming Huang, Zongyang Li, Dezhi Ran, Xinyu Wang, Ying Li, Assaf Marron, David Harel, Yuan Xie, and Tao Xie. 2025. DeCon: Detecting Incorrect Assertions via Postconditions Generated by a Large Language Model. *CoRR* abs/2501.02901 (2025). https://doi.org/10.48550/ARXIV.2501.02901 arXiv:2501.02901

[121] Zhiqiang Yuan, Mingwei Liu, Shiji Ding, Kaixin Wang, Yixuan Chen, Xin Peng, and Yiling Lou. 2024. Evaluating and Improving ChatGPT for Unit Test Generation. *Proc. ACM Softw. Eng.* 1, FSE (2024), 1703–1726. https://doi.org/10.1145/3660783

[122] Imam Nur Bani Yusuf and Lingxiao Jiang. 2024. Transducer Tuning: Efficient Model Adaptation for Software Tasks Using Code Property Graphs. *CoRR* abs/2412.13467 (2024). https://doi.org/10.48550/ARXIV.2412.13467 arXiv:2412.13467

[123] Zulfa Zakaria, Rodziah Binti Atan, Abdul Azim Abdul Ghani, and Nor Fazlida Mohd Sani. 2009. Unit Testing Approaches for BPEL: A Systematic Review. In *16th Asia-Pacific Software Engineering Conference, APSEC 2009, 1-3 December 2009, Batu Ferringhi, Penang, Malaysia*, Shahida Sulaiman and Noor Maizura Mohamad Noor (Eds.). IEEE Computer Society, 316–322. https://doi.org/10.1109/APSEC.2009.72

[124] He Zhang, Muhammad Ali Babar, and Paolo Tell. 2011. Identifying Relevant Studies in Software Engineering. *Inf. Softw. Technol.* 53, 6 (2011), 625–637. https://doi.org/10.1016/J.INFSOF.2010.12.010

[125] Jiyang Zhang, Yu Liu, Pengyu Nie, Junyi Jessy Li, and Milos Gligoric. 2024. Generating Exceptional Behavior Tests with Reasoning Augmented Large Language Models. *CoRR* abs/2405.14619 (2024). https://doi.org/10.48550/ARXIV.2405.14619 arXiv:2405.14619

[126] Quanjun Zhang, Chunrong Fang, Yuxiang Ma, Weisong Sun, and Zhenyu Chen. 2024. A Survey of Learning-based Automated Program Repair. *ACM Trans. Softw. Eng. Methodol.* 33, 2 (2024), 55:1–55:69. https://doi.org/10.1145/3631974

[127] Quanjun Zhang, Chunrong Fang, Weisong Sun, Yan Liu, Tieke He, Xiaodong Hao, and Zhenyu Chen. 2024. APPT: Boosting Automated Patch Correctness Prediction via Fine-Tuning Pre-Trained Models. *IEEE Trans. Software Eng.* 50, 3 (2024), 474–494. https://doi.org/10.1109/TSE.2024.3354969

[128] Quanjun Zhang, Chunrong Fang, Yang Xie, Yaxin Zhang, Yun Yang, Weisong Sun, Shengcheng Yu, and Zhenyu Chen. 2023. A Survey on Large Language Models for Software Engineering. *CoRR* abs/2312.15223 (2023). https://doi.org/10.48550/ARXIV.2312.15223 arXiv:2312.15223

[129] Quanjun Zhang, Chunrong Fang, Bowen Yu, Weisong Sun, Tongke Zhang, and Zhenyu Chen. 2024. Pre-Trained Model-Based Automated Software Vulnerability Repair: How Far are We? *IEEE Trans. Dependable Secur. Comput.* 21, 4 (2024), 2507–2525. https://doi.org/10.1109/TDSC.2023.3308897

[130] Quanjun Zhang, Chunrong Fang, Tongke Zhang, Bowen Yu, Weisong Sun, and Zhenyu Chen. 2023. Gamma: Revisiting Template-Based Automated Program Repair Via Mask Prediction. In *38th IEEE/ACM International Conference on Automated Software Engineering, ASE 2023, Luxembourg, September 11-15, 2023*. IEEE, 535–547. https://doi.org/10.1109/ASE56229.2023.00063

[131] Quanjun Zhang, Chunrong Fang, Yi Zheng, Ruixiang Qian, Shengcheng Yu, Yuan Zhao, Jianyi Zhou, Yun Yang, Tao Zheng, and Zhenyu Chen. 2025. Improving Retrieval-Augmented Deep Assertion Generation via Joint Training. *IEEE Transactions on Software Engineering* (2025), 1–15. https://doi.org/10.1109/TSE.2025.3545970

[132] Quanjun Zhang, Chunrong Fang, Yi Zheng, Yaxin Zhang, Yuan Zhao, Rubing Huang, Jianyi Zhou, Yun Yang, Tao Zheng, and Zhenyu Chen. 2025. Improving Deep Assertion Generation via Fine-Tuning Retrieval-Augmented Pre-trained Language Models. *ACM Trans. Softw. Eng. Methodol.* (Feb. 2025). https://doi.org/10.1145/3721128 Just Accepted.

[133] Quanjun Zhang, Ye Shang, Chunrong Fang, Siqi Gu, Jianyi Zhou, and Zhenyu Chen. 2024. TestBench: Evaluating Class-Level Test Case Generation Capability of Large Language Models. *CoRR* abs/2409.17561 (2024), arXiv–2409.

[134] Quanjun Zhang, Weifeng Sun, Chunrong Fang, Bowen Yu, Hongyan Li, Meng Yan, Jianyi Zhou, and Zhenyu Chen. 2024. Exploring Automated Assertion Generation via Large Language Models. *ACM Trans. Softw. Eng. Methodol.* (Oct. 2024). https://doi.org/10.1145/3699598 Just Accepted.

[135] Zhe Zhang, Xingyu Liu, Yuanzhang Lin, Xiang Gao, Hailong Sun, and Yuan Yuan. 2024. LLM-based Unit Test Generation via Property Retrieval. *CoRR* abs/2410.13542 (2024). https://doi.org/10.48550/ARXIV.2410.13542 arXiv:2410.13542

[136] Zibin Zheng, Kaiwen Ning, Qingyuan Zhong, Jiachi Chen, Wenqing Chen, Lianghong Guo, Weicheng Wang, and Yanlin Wang. 2025. Towards an Understanding of Large Language Models in Software Engineering Tasks. *Empir. Softw. Eng.* 30, 2 (2025), 50. https://doi.org/10.1007/S10664-024-10602-0

[137] Zhiyuan Zhong, Sinan Wang, Hailong Wang, Shaojin Wen, Hao Guan, Yida Tao, and Yepang Liu. 2024. Advancing Bug Detection in Fastjson2 with Large Language Models Driven Unit Test Generation. *CoRR* abs/2410.09414 (2024). https://doi.org/10.48550/ARXIV.2410.09414 arXiv:2410.09414

[138] Zhichao Zhou, Yutian Tang, Yun Lin, and Jingzhu He. 2024. An LLM-based Readability Measurement for Unit Tests' Context-aware Inputs. *CoRR* abs/2407.21369 (2024). https://doi.org/10.48550/ARXIV.2407.21369 arXiv:2407.21369