

Large Language Model-Based Agents for Software Engineering: A Survey

Junwei Liu, Kaixin Wang, Yixuan Chen, Xin Peng, Zhenpeng Chen, Lingming Zhang, Yiling Lou

Abstract—The recent advance in Large Language Models (LLMs) has shaped a new paradigm of AI agents, *i.e.*, LLM-based agents. Compared to standalone LLMs, LLM-based agents substantially extend the versatility and expertise of LLMs by enhancing LLMs with the capabilities of perceiving and utilizing external resources and tools. To date, LLM-based agents have been applied and shown remarkable effectiveness in Software Engineering (SE). The synergy between multiple agents and human interaction brings further promise in tackling complex real-world SE problems. In this work, we present a comprehensive and systematic survey on LLM-based agents for SE. We collect 106 papers and categorize them from two perspectives, *i.e.*, the SE and agent perspectives. In addition, we discuss open challenges and future directions in this critical domain. The repository of this survey is at <https://github.com/FudanSELab/Agent4SE-Paper-List>.

Index Terms—Large Language Model, AI Agent, Software Engineering



1 INTRODUCTION

Large Language Models (LLMs) [1] have achieved remarkable progress and demonstrated potential of human-like intelligence. In recent years, LLMs have been widely applied in Software Engineering (SE). As shown by recent surveys [2], [3], LLMs have been adopted and shown promising performance in various software development and maintenance tasks, such as program generation [4]–[8], software testing [9]–[11] and debugging [12]–[17], and program improvement [18]–[20].

AI Agents are artificial entities that can autonomously perceive and act on surrounding environments so as to achieve specific goals [21]. The concept of AI agents has been evolving for a long time (*e.g.*, early agents are constructed on symbolic logic or reinforcement learning [22]–[25]). Recently, the remarkable progress in LLMs has further shaped a new paradigm of AI agents, *i.e.*, LLM-based agents, which leverage LLMs as the central agent controller. Different from standalone LLMs, LLM-based agents extend the versatility and expertise of LLMs by equipping LLMs with the capabilities of perceiving and utilizing external resources and tools, which can tackle more complex real-world goals via collaboration between multiple agents or involvement of human interaction.

In this work, we present a comprehensive and systematic survey on LLM-based agents for SE. We collect 106 papers and categorize them from two perspectives, *i.e.*, both the SE and agent perspectives. Additionally, we discuss the open challenges and future directions in this domain.

- J. Liu, K. Wang, Y. Chen, X. Peng, and Y. Lou are with the Department of Computer Science, Fudan University, China. E-mails: {juliu24, kxwang23, yixuanchen23}@m.fudan.edu.cn, {pengxin, yilinglou}@fudan.edu.cn
- Z. Chen is with the School of Computer Science and Engineering, Nanyang Technological University, Singapore. E-mail: zhenpeng.chen@ntu.edu.sg
- L. Zhang is with the Department of Computer Science, University of Illinois Urbana-Champaign, USA. E-mail: lingming@illinois.edu

From the *SE* perspective, we analyze how LLM-based agents are applied across different software development and improvement activities, including individual tasks (*e.g.*, requirements engineering, code generation, static code checking, testing, and debugging) as well as the end-to-end procedure of software development and improvement. From this perspective, we provide a comprehensive overview of how SE tasks are tackled by LLM-based agents.

From the *agent* perspective, we focus on the design of components in LLM-based agents for SE. Specifically, we analyze key components, including planning, memory, perception, and action, in these agents. Beyond basic agent construction, we also analyze multi-agent systems, including their agent roles, collaboration mechanisms, and human-agent collaboration. From this perspective, we summarize the characteristics of different components of LLM-based agents when applied to the SE domain.

In summary, this survey makes the following contributions:

- It provides the first comprehensive survey of 106 papers that apply LLM-based agents to SE.
- It analyzes how existing LLM-based agents are designed and applied for software development and maintenance from both the SE and agent perspectives.
- It discusses research opportunities and future directions in this critical domain.

Survey Structure. Figure 1 summarizes the structure of this survey. Section 2 introduces background knowledge, while Section 3 presents the methodology. Section 4 and Section 5 present the relevant work from the SE perspective and the agent perspective, respectively. Finally, Section 6 discusses the potential research opportunities.

2 BACKGROUND AND PRELIMINARY

In this section, we first introduce the background about the basic and advanced LLM-based agents, and then we discuss the related surveys.

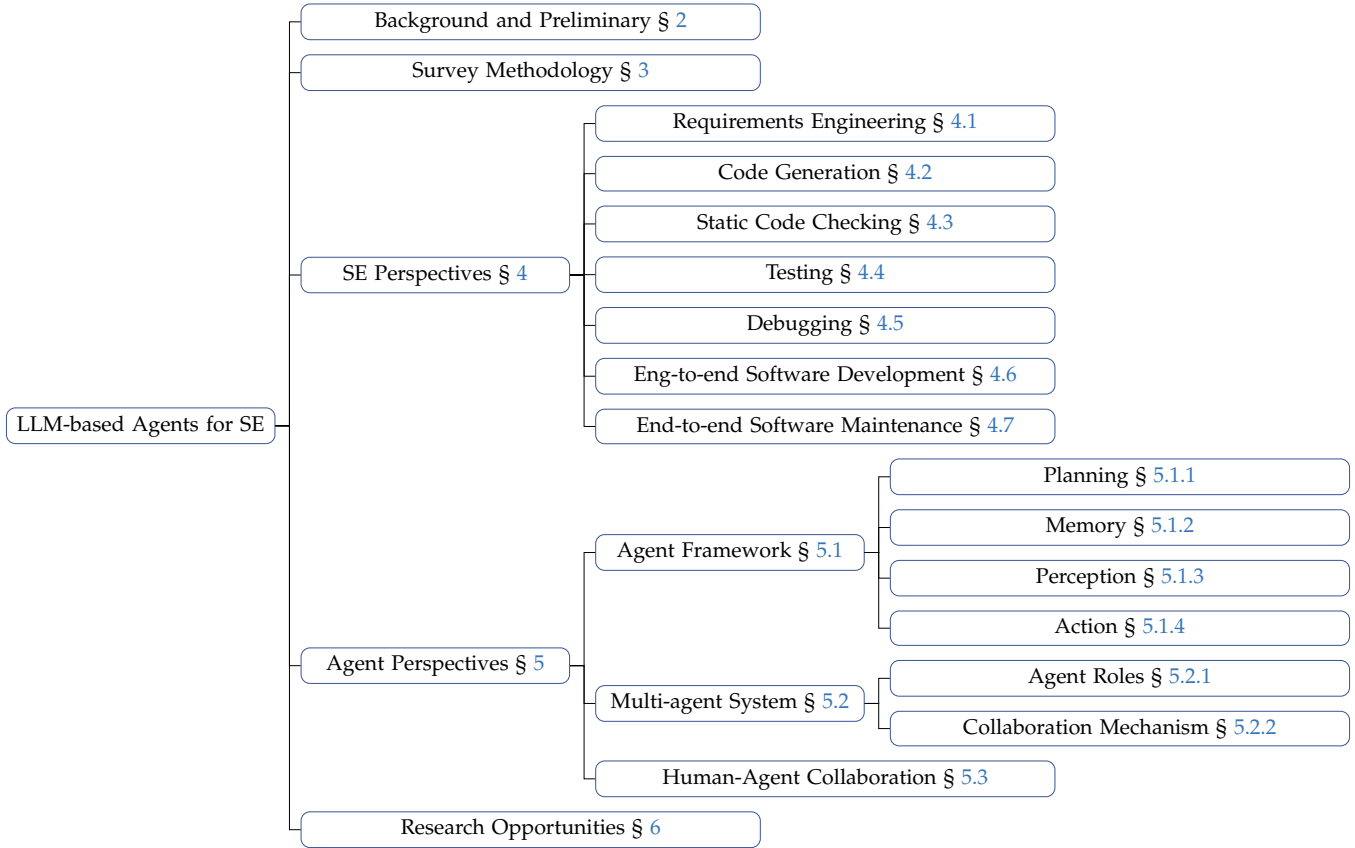


Fig. 1: Structure of This Survey

2.1 Basic Framework of LLM-based Agents

LLM-based agents are typically composed of four key components: *planning*, *memory*, *perception*, and *action* [21]. The planning and memory serve as the key components of the *LLM-controlled brain*, which interacts with the environment through the perception and action components to achieve specific goals. Figure 2 illustrates the basic framework of LLM-based agents.

Planning. The planning component decomposes complex tasks into multiple sub-tasks and schedules the sub-tasks to achieve final goals. In particular, agents can (i) generate a plan without adjustment by different reasoning strategies, or (ii) adjust a generated plan with the external feedback (*e.g.*, environmental feedback or human feedback).

Memory. The memory component records the historical thoughts, actions, and environmental observations generated during the agent execution [21], [26], [27]. Based on the accumulated memory, agents can revisit and utilize the previous records and experience, so as to tackle the complex tasks more effectively. The memory management (*i.e.*, how to represent the memory) and utilization (*i.e.*, how to read/write or retrieve the memory) are essential, which directly impact the efficiency and effectiveness of the agent system.

Perception. The perception component receives the information from the environment, which can facilitate better planning. In particular, agents can perceive multi-modal inputs, *e.g.*, textual inputs, visual inputs, and auditory inputs.

Action. Based on the planning and decisions made by

the brain, the action component conducts concrete actions to interact with and impact the environment. One essential mechanism in action is to control and utilize external tools, which can extend the inherent capabilities of LLMs by accessing more external resources and extending the action space beyond textual-alone interaction.

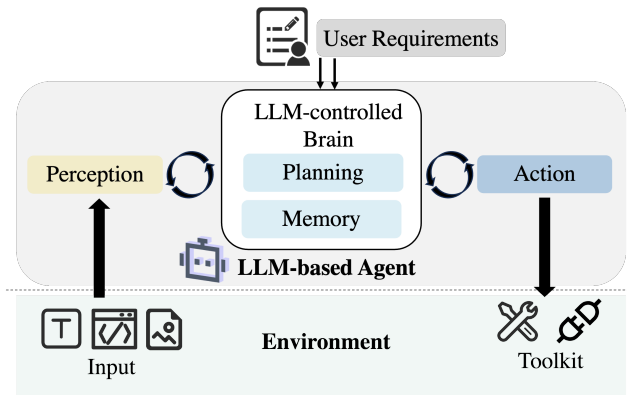


Fig. 2: Basic Framework of LLM-based Agents

2.2 Advanced LLM-based Agent Systems

Multi-agent Systems. While a single-agent system can be specialized to solve one certain task, enabling the collaboration between multiple agents (*i.e.*, *multi-agent systems*) can

further solve more complex tasks associated with diverse knowledge domains. In particular, in a multi-agent system, each agent is assigned a distinct role and relevant expertise, making it specifically responsible for different tasks; in addition, the agents can communicate with each other and share the progress/information as the task proceeds. Typically, agents can work collaboratively (i.e., by working on different sub-tasks to achieve a final goal) or competitively (i.e., by working on the same task while debating adversarially).

Human-Agent Coordination. Agent systems can further incorporate the instructions from humans and then proceed with tasks under human guidance. This human-agent coordination paradigm facilitates better alignment with human preference and uses human expertise. In particular, during human-agent interaction, humans can not only provide agents with task requirements and feedback on the current task status, but also cooperate with agents to achieve goals together.

2.3 Related Surveys

LLM-based agents in general domains have been widely discussed and surveyed [21], [26], [28]–[32]. Different from these surveys, this survey focuses on the design and application of LLM-based agents specifically for the software engineering domain. In software engineering domain, there have been several surveys or literature reviews on the general application of LLMs in software engineering [2], [3], [10], [32], [33]. Different from these surveys, this survey specifically focuses on the agent perspective and is more comprehensive on the application of LLM-based agents for software engineering. In addition, He *et al.* [34] present a vision paper on the potential applications and emerging challenges of multi-agent systems for software engineering. Different from the vision paper, this work focuses on conducting a comprehensive survey of existing agent systems (including both single agent and multi-agent systems). In summary, to the best of our knowledge, this is the first survey specifically focusing on the literature on LLM-based agents for software engineering.

3 SURVEY METHODOLOGY

This section defines the scope of the survey and describes our approach to collecting and analyzing papers within the scope.

3.1 Survey Scope

We focus on the papers that apply *LLM-based agents* to tackle *SE tasks*. In the following, we specify the terms.

- *SE tasks.* Following previous surveys on the application of LLMs in SE [2], [3], we focus on all SE tasks along the software life cycle, including requirements engineering, software design, code generation, software quality assurance (i.e., static checking and testing), and software improvement.
- *LLM-based agents.* A standalone LLM can work as a naive “agent” since it can take textual inputs and produce textual outputs, leaving it no clear boundary between LLMs and LLM-based agents. However, this could result

in an overly broad scope and significant overlap with existing surveys on LLM applications in SE [2], [3]. Based on the widely-adopted consensus about AI agents, the key characteristic of agents is their ability to autonomously and iteratively perceive feedback from, and act upon, a dynamic environment [21]. To ensure a more focused discussion from the perspective of agents, this survey focuses on LLM-based agents that not only incorporate LLMs as the core of their “brains”, but also have the capacity to iteratively interact with the environment, taking feedback and acting in real time.

More specifically, we apply the following inclusion and exclusion criteria for paper collection.

- **Inclusion criteria.** A paper will be included in our survey if it meets any of the following criteria: (i) The paper proposes a technique, framework, or tool for addressing specific SE tasks using LLM-based agents; (ii) The paper presents a general technique, framework, or tool applicable across various domains, provided that its evaluation includes at least one SE task; (iii) The paper presents an empirical study evaluating LLM-based agents on specific SE tasks.
- **Exclusion criteria.** A paper will be excluded from our survey if it meets any of the following criteria: (i) The paper does not involve any SE tasks; (ii) The paper only discusses LLM-based agents in the context of discussion or future work, without integrating them into the main approach; (iii) The paper only uses a standalone LLM for processing textual inputs and generating textual outputs, without any iterative interaction with the environment.

3.2 Paper Collection

Our paper collection process includes two steps: keyword searching and snowballing.

3.2.1 Keyword Searching

We follow established practices in SE surveys [35]–[39] by using the DBLP database [40] for paper collection. Recent research [39] has demonstrated that papers gathered from other prominent publication databases are typically a subset of those available on DBLP, which encompasses over 7 million publications from more than 6,500 academic conferences and 1,850 journals in computer science [41]. DBLP also covers arXiv [42], a widely adopted open-access repository.

We employ an iterative trial-and-error approach, which is widely adopted in SE surveys [35], [43], to determine search keywords. Initially, all authors, with relevant research experience/publication in LLM and SE, convene to suggest papers relevant to our scope, yielding an initial set of relevant papers. Subsequently, the first two authors review the titles, abstracts, and introductions of these papers to identify additional keywords. We then conduct brainstorming sessions to expand and refine our search strings, incorporating related terms, synonyms, and variations. This process enables iterative enhancement of our search keyword list. The final keywords include (“agent” OR “llm” OR “language model”) AND (“api” OR “bug” OR “code” OR “coding” OR “debug” OR “defect” OR “deploy” OR “evolution” OR “fault” OR “fix” OR “maintenance” OR “program” OR “refactor”

TABLE 1: Statistics of Paper Collection

Keyword			Hits
agent	llm	language model + api	83
agent	llm	language model + bug	98
agent	llm	language model + code	915
agent	llm	language model + coding	70
agent	llm	language model + debug	95
agent	llm	language model + defect	22
agent	llm	language model + deploy	295
agent	llm	language model + evolution	1,349
agent	llm	language model + fault	685
agent	llm	language model + fix	318
agent	llm	language model + maintenance	64
agent	llm	language model + program	1,969
agent	llm	language model + refactor	15
agent	llm	language model + repair	137
agent	llm	language model + requirement	451
agent	llm	language model + software	2,151
agent	llm	language model + test	976
agent	llm	language model + verification	525
agent	llm	language model + vulnerab	144
After manual inspection			67
After snowballing			106

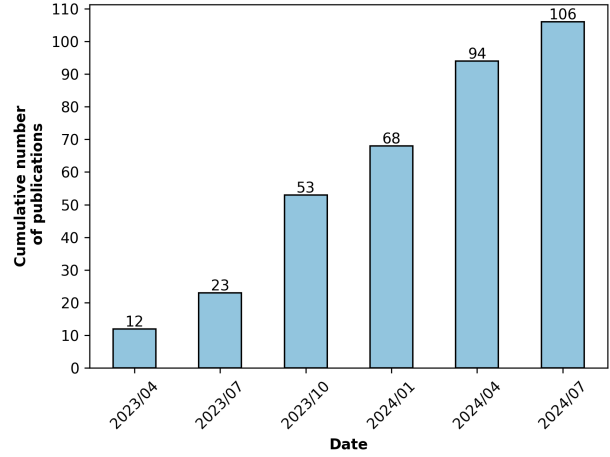


Fig. 3: Cumulative Number of Papers Over Time

OR “repair” OR “requirement” OR “software” OR “test” OR “verification” OR “vulnerab”).

Based on the keywords, we conduct 57 searches on DBLP on July 1st, 2024, and obtain 10,362 hits. Table 1 presents the statistics of papers collected through keyword searching. The first two authors manually review each paper to filter out those not within the scope of this survey. As a result, we identify 67 relevant papers through this process.

3.2.2 Snowballing

To enhance the comprehensiveness of our survey, we adopt snowballing approaches to identify papers that are transitively relevant and expand our paper collection [35]. Specifically, between July 1 and July 10, 2024, we conduct both backward and forward snowballing. Backward snowballing involves examining references in each collected paper to identify relevant ones within our scope, while forward snowballing uses Google Scholar to find relevant papers citing the collected ones. This iterative process continues until no new relevant papers are found. In this process, we retrieve an additional 39 papers.

3.3 Statistics of Collected Papers

As shown in Table 1, we have collected a total of 106 papers for this survey. Figure 3 presents the cumulative number of papers published over time, up to July 10, 2024. We observe that there is a continuous increase of research interest in this field, highlighting the necessity and relevance of this survey. Additionally, Figure 4 shows the distribution of publication venues for the papers, covering diverse research communities such as software engineering, artificial intelligence, and human-computer interaction. In particular, the majority of the papers are from arXiv and have not yet undergone peer review. This is expected, as this field is emerging and still undergoing rapid development.

4 ANALYSIS FROM SE PERSPECTIVES

In this section, we organize the collected papers from the perspective of different SE tasks. Figure 5 presents the SE

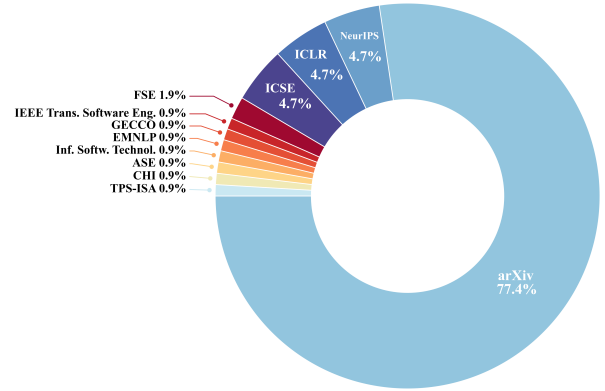


Fig. 4: Distribution of Publication Venues of All Papers

tasks along the common life cycle of software development and maintenance.

It is worth noting that, LLM-based agents can be designed not only to tackle individual SE tasks but also to support end-to-end software development or maintenance processes involving multiple SE activities. From the collected papers, we observe LLM-based agents designed for (i) *end-to-end software development* and (ii) *end-to-end software maintenance*. Specifically, agents for end-to-end software development can generate a complete program based on requirements by performing multiple SE tasks, such as requirements engineering, design, code generation, and code quality assurance (e.g., verification, static checking, and testing); agents for end-to-end software maintenance can generate patches for user-reported issues by supporting multiple SE maintenance activities, such as debugging (e.g., fault localization and repair) and feature maintenance. As shown in previous papers [2], [3], standalone LLMs are primarily specialized in tackling single SE tasks and are generally inadequate for complex end-to-end software development and maintenance processes. In contrast, LLM-based agents, through their components (i.e., planning, memory, perception, and action), coordination among multiple agents, and human interaction, provide the autonomy

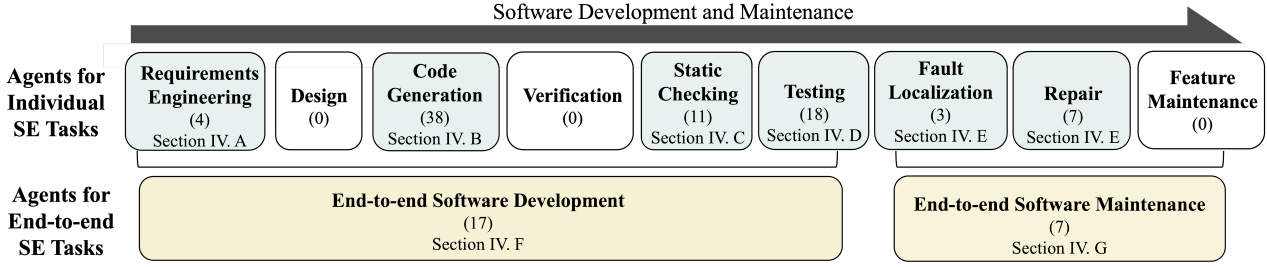


Fig. 5: Agent Distribution along Software Development and Maintenance Tasks

and flexibility necessary to tackle these complex tasks.

Distribution of LLM-based agents in different SE activities. In Figure 5, the numbers in brackets indicate the count of collected papers in each category. Notably, if LLM-based agents are designed for end-to-end software development or maintenance, they are only reported at the end-to-end level rather than at the level of individual tasks. Overall, we observe that the majority of LLM-based agents focus on individual-level SE tasks, especially for code generation and code quality assurance (*e.g.*, static checking and testing); in addition, a portion of agents are designed for end-to-end software development or maintenance tasks, indicating the promise of LLM-based agents in tackling more complex real-world SE tasks.

4.1 Requirements Engineering

Requirements Engineering (RE) is a crucial phase for initializing the software development procedure. Generally, it can cover the following phases [44]–[46].

- *Elicitation*: New requirements are elicited and collected.
- *Modeling*: Abstract yet interpretable models are used to describe requirements, *e.g.*, Unified Modeling Language (UML) [47] and Entity-Relationship-Attribute (ERA) model [48].
- *Negotiation*: Negotiation plays a crucial role in facilitating communication of different stakeholders and ensuring consistency, especially in conflicting requirements.
- *Specification*: Requirements are determined and documented in a formal format.
- *Verification*: Requirements and models are validated to ensure they fully and unambiguously reflect the intent of stakeholders.
- *Evolution*: Requirements evolution refers to the ongoing process of refining and adapting requirements in response to changing needs and conditions.

In real-world software development, RE can take lots of manual efforts due to the strong demand for massive interactions with various stakeholders. Although researchers have leveraged deep learning models (including standalone LLMs) to boost requirements engineering activities, most of them still remain on individual tasks during RE, such as classification [49], specification [50], information retrieval [51], evaluation [52], and enhancement [53] of existing requirements. Recently, multi-agent systems are designed to automate individual phases or multiple phases. Table 2 summarizes existing LLM-based agents specifically designed for RE, and Figure 6 illustrates their common pipeline.

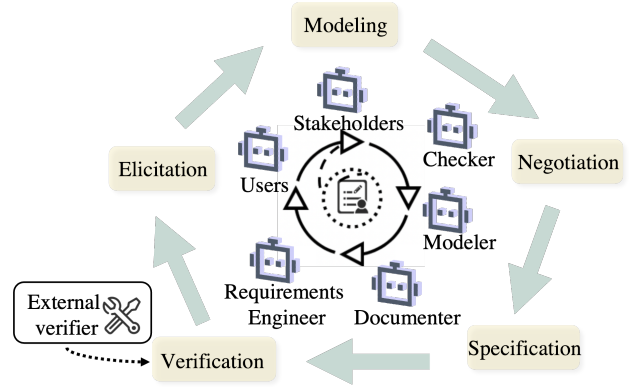


Fig. 6: Pipeline of LLM-based Agents for Requirements Engineering

A. LLM-based Agents for Individual RE Phases. Elicitation [54] is a multi-agent framework for the requirement elicitation phase, which aims at mining requirements as completely as possible. Elicitation first initializes multiple agents with different personas within the designed context to cover various user viewpoints, and then makes them simulate interactions with the target product while documenting records (*i.e.*, action, observation, and challenge) in each step. The potential requirements are eventually identified through the agent interviews and filtered on the provided criteria. Experimental results indicate that Elicitation can uncover and categorize hidden needs while reducing costs compared with conventional methodologies such as user studies.

SpecGen [55] is a system designed for generating requirement specifications of a given program. It contains two stages: conversation-driven specification generation and mutation-based specification generation. In the first stage, an agent generates Java Modeling Language (JML) requirement specifications [58], which are further validated by the external OpenJML verifier [59]. The feedback is integrated into the prompt and kicks off iterative generation. In the second stage, specifications that fail the validation are mutated and verified to generate more diverse specifications. Experimental results demonstrate that SpecGen outperforms state-of-the-art approaches.

B. LLM-based Agents for Multiple RE Phases. In addition to automating individual RE phases, some agents address multiple RE phases. For example, Arora *et al.* [56]

TABLE 2: Existing LLM-based Agents for Requirements Engineering

Agents	Multi-Agent	Covered RE Phases					
		Elicitation	Modeling	Negotiation	Specification	Verification	Evolution
Elictron [54]	✓	✓					
SpecGen [55]	×				✓		
Arora <i>et al.</i> [56]	✓	✓		✓	✓	✓	
MARE [57]	✓	✓	✓		✓	✓	

propose a multi-agent system to cover four phases of RE: elicitation, specification, analysis (synonymous with negotiation), and validation. For each phase, they explore the roles that can be simulated by LLM-based agents and systematically analyze the strengths and weaknesses. In the elicitation phase, agents play the role of stakeholders or requirement engineers and collect preliminary requirements, which are then formatted into structured documentation by another agent in the specification phase. Then, the analysis phase involves multiple stakeholder agents, aiming to evaluate, prioritize, and refine requirements. The validation phase marks the final validation by the stakeholder agents, and then the requirement documentation is decided and ready for subsequent design/implementation.

MARE [57] is another multi-agent framework that covers multiple RE phases, including elicitation, modeling, verification, and specification. In the elicitation phase, a set of stakeholder agents express their needs, which would then be organized into a draft by the collector agent. Subsequently, the modeler agent identifies entities and relationships in the draft and constructs a requirement model; In the verification phase, the checker agent assesses the quality of the current requirements draft on its criteria and hands it over to the documenter agent, which will write the requirement specifications or report errors. All of these agents are equipped with predefined actions and can communicate within a shared workspace, enabling the seamless exchange of intermediate information.

4.2 Code Generation

Code generation has been extensively explored with the development of AI technology [33]. Due to being pre-trained on massive textual data (especially large code corpus), LLMs demonstrate promising effectiveness in generating code for given code contexts or natural language descriptions. Nevertheless, the code generated by LLMs can sometimes be unsatisfactory due to issues such as the notorious hallucination [60]. Therefore, beyond simply leveraging standalone LLMs for code generation, researchers also build LLM-based agents that can enhance the capabilities of LLMs via planning and iterative refinement. Figure 7 illustrates how existing LLMs extend standalone LLMs in code generation.

4.2.1 Code Generation with Planning

LLM-based agents employ advanced planning methods to extend the code generation capabilities of LLMs. Chain-of-thought (CoT) [99] is the most popular strategy, which decomposes the code generation task into sub-tasks and achieves higher generation correctness [61], [70], [76], [84], [86], [91]–[93], [95]. For example, CodeCoT [84] leverages

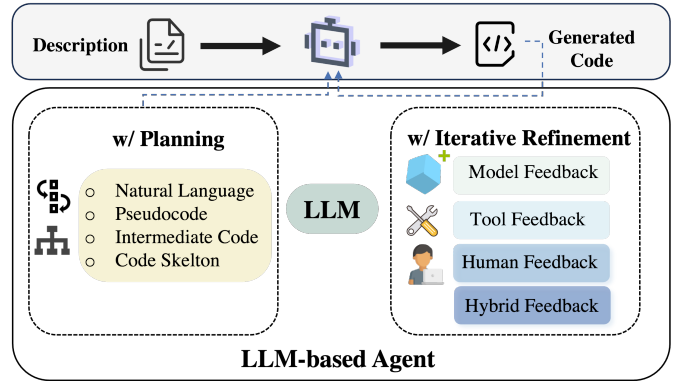


Fig. 7: Pipeline of LLM-based Agents for Code Generation

CoT to break down the given requirements into steps that are described in natural language and then convert them to code. Some works employ dynamic planning strategies, where observations of the current action determine the next step [62], [85], [87], [88]. For example, CodePlan [88] employs an adaptive planning algorithm that dynamically detects the affected code snippets in the repository and adapts the plan. However, there is still a gap between the narrative-based steps and the final generated code. Therefore, some works propose different representations to describe the planning steps, including *pseudocode* [95], *intermediate code* [61], or *code skeleton* [93], [97]. For example, AgentCoder [95] prompts the agent to generate pseudocode after problem understanding and algorithm selection phases, which serves as a draft for the final code. In addition, some existing works explore multi-path planning strategies. For example, LATS [76] simulates all possible generation paths as a tree and optimizes the plan with the Monte Carlo Tree Search algorithm. In MapCoder [98], the planning agent generates multiple plans along with confidence scores for sorting. The highest-scoring plan is used to generate the target code. If the code is erroneous, the plan with the next highest confidence is selected to continue the iterative generation process.

4.2.2 Code Generation with Iterative Refinement

One essential capability of agents is to act on the feedback from the environment. In the code generation scenario, some agents dynamically refine the previously-generated code based on the feedback via multiple iterations. We organize the relevant research based on the feedback sources, including model feedback, tool feedback, human feedback, and hybrid feedback. Table 3 summarizes existing LLM-based agents for code generation with the iterative refinement.

A: Model Feedback. Model feedback can be classified into peer-reflection and self-reflection.

TABLE 3: Existing LLM-based Agents for Code Generation

Agents	Multi-Agent	Iterative Refinement			
		Model Feedback	Tool Feedback	Human Feedback	Hybrid Feedback
Parsel [61]	✓				
Reflexion [62], Self-Repair [63], AutoGen [64], INTERVENOR [65], TGen [66], AutoCoder [67]	✓	✓	✓		✓
CAMEL [68], Li <i>et al.</i> [69], DyLAN [70]	✓	✓			
SELF-DEBUGGING [71], SEIDR [72], μFiX [73], AlphaCodium [74], LDB [75], LATS [76], RRR [77]	×	✓	✓		✓
ToolCoder [78], SELF-EVOLVE [79], KPC [80], LEMUR [81], CODEAGENT [82], LLM4TDD [83], CodeCoT [84], CodeAct [85], CoCoST [86], InterCode [87], CodePlan [88], TOOLGEN [89]	×		✓		
Self-Refine [90]	×	✓			
Flows [91]	✓	✓	✓	✓	✓
MINT [92]	×	✓	✓	✓	
CodeChain [93]	×	✓	✓		
ClarifyGPT [94]	×		✓	✓	
AgentCoder [95], Gentopia [96], SoA [97], MapCoder [98]	✓		✓		

A.1: Peer-reflection. Peer-reflection refers to information exchange and interaction between multiple models. The most common paradigm is collaboration in multi-agent systems through role specialization and structured communication (e.g., code review) [62], [64], [68]. This approach underscores the specialized responsibilities of each role and how they exchange information based on these responsibilities. Besides, when generating initial code, some works produce multiple results at once [69]. Thus, a selection mechanism is employed to retain the most suitable result. Li *et al.* [69] use Bilingual Evaluation Understudy (BLEU) to calculate and aggregate the similarity scores of each initial code with the rest, retaining the result with the highest score. Moreover, there is also a modality that involves treating each agent role equally in expressing their opinions or engaging in debates to solve problems. DyLAN [70] allows multiple agents to engage in dynamic interactions over multiple rounds, organized into a multi-layer feed-forward network. The network employs an additional LLM ranker to analyze the responses of agents from the previous layer and selects the best-performing agent to continue in subsequent interactions.

A.2: Self-reflection. Apart from the interaction between models, there are works that conduct self-refinement of a single model [71], [73], [90]. This means that the current modification is based on the previous output of the model, iteratively optimizing through this approach. Le *et al.* [93] guide LLMs to generate modularized code, leveraging cluster representatives from previously generated sub-modules in each iteration. SELF-DEBUGGING [71] draws inspiration from the rubber duck debugging method used by programmers. During the explanation phase, the model provides a line-by-line explanation of the generated initial code. As Wang *et al.* [92] mention in their work, all models benefit from natural language feedback, with absolute performance gains by 2–17% each additional turn of natural language feedback.

B: Tool Feedback. The code generated by models can

be of limited quality with numerous uncertainties. One solution to address this challenge is to equip LLM-based agents with tools that can collect informative feedback and assist the agents to generate and refine code.

B.1: Dynamic Execution Tools. One common group is to invoke the compiler, interpreter, and execution engine to directly compile or execute the code. This approach leverages the outputs and run-time behaviors, such as test results or compilation errors, as feedback for code improvement [79], [81]–[87], [92]–[98].

B.2: Static Checking Tools. Agents can get more restricted knowledge on code constraints by applying code analysis tools. For example, some agents apply static analysis tools to obtain syntactically-valid program symbols/tokens [89] or dependencies between code during code generation [82], [88]. Including the analyzed information into the prompt can guide LLMs towards generating valid code.

B.3: Retrieval Tools. Agents can get the access to rich external resources by applying retrieval or searching tools. For example, some agents retrieve local knowledge repositories [86] such as private API documentations [78], [82] to facilitate better code generation; in addition, some apply online search engines [78], [82], [86], [96] or web crawling [80] to collect information such as content from relevant websites (e.g., *StackOverflow* and *datagy.io*) [78], [82], [86], [96] and official online documentations [80], [86]. Including the retrieved resources into the prompt can provide additional knowledge for language models. Notably, ToolCoder [78] integrates the agent with online search and local documentation search tools that provide helpful information for both public and private APIs, alleviating the hallucination of LLMs.

C: Human Feedback. Another approach involves incorporating human feedback into the process, as humans play a critical role in clarifying ambiguous requirements. For instance, in software development, humans can check whether the generated code aligns with their initial intent.

Any discrepancies are often attributed to vagueness or incompleteness in the requirements, prompting a revision of the requirement documents [91], [92]. To further minimize human involvement, some methods enable the agent to handle the task of observing execution results. For example, ClarifyGPT [94] automatically identifies potential ambiguities in the manually-given requirements and proactively poses relevant questions for humans; then the responses from humans are further used to refine the requirements.

D: Hybrid Feedback. Agents can also incorporate multiple types of feedback and progressively enhance each other as hybrid feedback. For example, a common approach is to combine tool feedback and model feedback. Specifically, an LLM receives error messages returned after executing a program or test case, and utilizes its contextual understanding to provide corresponding feedback output (e.g., explanations, suggestions, instructions, etc.) [62]–[67], [71]–[77], [91], [97]. For example, in the multi-agent system INTERVENOR [65], a teacher coder is designated to observe the program execution results and provide error explanations and bug-fixing plans for the student coder to understand and regenerate the code. Furthermore, to accurately pinpoint issues, some works provide more fine-grained environment feedback for the subsequent processing of the model. LDB [75] constructs a control flow graph, which divides the program into multiple blocks. It uses breakpoints to obtain the runtime values of variables, and the model compares these values against the requirements, assessing each block for anomalies.

4.3 Static Code Checking

Static code checking refers to examining the quality of code without executing the code. In particular, static code checking has been essential in the modern continuous integration pipeline, as it is efficient to identify diverse categories of code quality issues (e.g., different bugs, vulnerabilities, or code smells) before extensively executing the tests. In practice, it is common to adopt static analysis techniques to automatically detect bugs/vulnerabilities (i.e., static bug detection) or involve peer reviews to check the quality of code (i.e., code review).

4.3.1 Static Bug Detection

Preliminary studies [2], [3] show that LLMs can help identify potential quality issues in the given code under inspection. For example, fine-tuning LLMs on existing buggy/correct code or simply prompting LLMs has demonstrated promising effectiveness in identifying the bugs, vulnerabilities, or code smells in the given code snippets [15], [100]. However, given the diversity and complexity of the root causes of different code issues as well as the long code contexts under inspection, standalone LLMs exhibit limited accuracy and recall in the real-world static code checking scenario [101]. Recently, researchers have built LLM-based agents to enhance the capabilities of standalone LLMs in bug or vulnerability detection. Table 4 summarizes these agents and Figure 8 illustrates their common pipeline.

A. Co-inspection with Multi-agent. One effective vulnerability detection strategy focuses on the perspective of

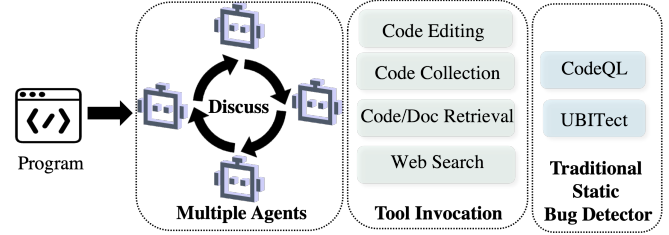


Fig. 8: Pipeline of LLM-based Agents for Static Bug Detection

multi-agent collaboration. Mao *et al.* [111] propose an approach for vulnerability detection through mutual discussion and consensus among different LLMs representing the testers and developers respectively. It mimics the real-world code review process that involves the collaboration of various roles within a team, enhancing the detection effectiveness through the interaction and reflection between LLMs. Moreover, GPTLENS [104] is an adversarial yet synergic framework for detecting vulnerabilities in smart contracts. It is designed as a two-stage method involving several auditor agents and a critic agent, all adopting GPT-4 as the backend. The auditor agents randomly generate potential vulnerabilities and corresponding reasoning as thoroughly as possible, while the critic scrutinizes and scores the candidates based on specific criteria. Their evaluation on 13 real-world smart contract CVEs indicates the effectiveness by an improvement of up to 76.9% in the identification rate. Fan *et al.* [105] conceive the Intelligent Code Analysis Agent (ICAA) concept for static code analysis. ICAA is an integration of AI models (e.g., the LLMs), engineering process designs, and traditional non-AI components (e.g., static analysis tools). As a dynamic decision-making system, ICAA can be composed of multiple sub-agents to enhance its functionality. Two examples of ICAA implementations include bug detection and code-intention consistency checking, both of which rely on the collaboration of various sub-agents, such as the ReAct-based Analysis Agent and the Report Agent.

B. Additional Knowledge from Tool Execution. Other approaches concentrate on strengthening the capabilities of LLMs via tool invocation. ART [102] is a general framework that boosts LLMs for unseen tasks with multi-step planning and effective tool utilization. ART includes a tool library (e.g., search tools), which can facilitate task decomposition and the appropriate invocation of tools. In addition, ART accommodates the integration of human inputs, allowing for seamless updates to its libraries. ART is evaluated on the multiple downstream tasks such as bug detection, and outperforms both few-shot prompting and the automated generation of CoT reasoning. Sun *et al.* [110] propose LLM4Vuln to enhance the vulnerability reasoning capabilities of LLMs by decoupling and augmenting them. The agent improves in several aspects. It retrieves external knowledge from both the raw text of vulnerability reports and summarized key sentences, which include information about the functionality of the vulnerable code and its root cause. This information is stored in a vector database. Additionally, the agent invokes tools to actively seek further context about the

TABLE 4: Existing LLM-based Agents for Static Bug Detection

Agents	Multi-Agent	Tool Utilization		Dataset	Target Program	Bug Category
		Tool Category	Specific Tools			
ART [102]	×	Custom Toolkit	Code Generation Tool Code Execution Tool	BigBench [103]	Python Program	Code Errors
GPTLENS [104]	✓	-	-	Self-curated	Smart Contract	Smart Contract Vulnerability
ICAA [105]	✓	Custom Toolkit	Context Splitting Tool Code Retrieval Tool Document Retrieval Tool Web Search Tool	NFBugs [106] Self-curated	Python Program Java Program	Non-functional Bugs API Misusage
E&V [107]	×	Static Analysis	Clang [108]	Sampled syzbot [109]	Linux Kernel	Kernel Address Sanitizer Bugs
LLM4Vuln [110]	×	Custom Toolkit	Database Retrieval Tool Context Collection Tool	Self-curated	Smart Contract	Smart Contract Vulnerability
Mao <i>et al.</i> [111]	✓	-	-	SySeVR [112]	C/C++ Program	Library/API Function Call Arithmetic Expression Array Usage Pointer Usage
IRIS [113]	×	Static Analysis	CodeQL [114]	CWE-Bench-Java [113]	Java Program	Path-Traversal OS Command Injection Cross-Site Scripting Code Injection
LLIFT [115]	×	Static Analysis	UBITect [116]	Rnd-300 [115]	Linux Kernel C Program	UBI Bugs

target code (e.g., function or variable definitions) through function calls.

C. Combined with Traditional Static Bug Detection. Some researchers have combined LLM-based agents with traditional static checking techniques to improve their static bug detection capability. For example, LLIFT [115] is an LLM-assisted Use-Before-Initialization (UBI) bug detection tool. Based on undecided bugs reported by the powerful static analysis tool UBITect, LLIFT further leverages the capability of LLMs in code comprehension and summarization to identify UBI bugs in the Linux kernel. However, LLMs have the inherent limitation in accepting and understanding long input context (e.g., numerous functions in the Linux kernel) as well as the hallucination and stochasticity issues. LLIFT addresses these issues by augmenting the basic LLM with some agentic techniques. For example, instead of flooding the LLMs with all possible related code, it only responds to the LLM’s demand for specific function definitions through static analysis, achieving a balance between context length and information completeness. Besides, it employs in-context learning, task decomposition, self-validation, and majority voting strategies to further alleviate hallucination and stochasticity and guarantee detection accuracy. This framework identifies 13 UBI bugs from 1,000 potential UBI bug instances reported by UBITect [116], with a precision rate of 50%. E&V [107] is an agent designed for conducting static analysis of code in the Linux kernel. The high-level workflow of E&V is a loop of employing an LLM-based agent for static analysis through pseudo-code execution, verifying the output of pseudo-code, and providing feedback for re-analysis. To avoid fact hallucination caused by missing necessary code snippets (e.g., caller and callee functions in inter-procedural analysis), a source code retrieval component is introduced to fetch needed functions or structures through traditional static analysis tools (e.g., Clang [108]). E&V has been evaluated against 170 Linux Kernel bugs and correctly pinpointed the

blamed function in 81.2% of the cases. IRIS [113] is an agent augmented with CodeQL (a static analysis tool) [114] for vulnerability detection. IRIS first utilizes CodeQL to extract candidate APIs in the given repository. Then, it labels these APIs as potential sources or sinks of the given vulnerability via querying the LLM-based agent, which will be further handed over to CodeQL for detecting vulnerable paths. The final verdict is achieved by prompting the LLM agent to analyze the vulnerable paths and surrounding code of the source and sink.

4.3.2 Code Review

Developers review each other’s code changes to ensure and improve the code quality before merging the changes into the branch. To mitigate the manual efforts in code review, researchers leverage learning approaches to automate the code review procedure. In particular, code review is formulated as a binary classification problem (i.e., code quality classification [117]) or a sequence-to-sequence generation problem (i.e., review comment generation [118]), which are tackled by fine-tuning or prompting deep learning models (including LLMs). Different from these works, LLM-based agents mimic the real-world peer review procedure by including multiple agents as different code reviewers. Table 5 summarizes existing agents for code review.

CodeAgent [119] is a multi-agent system that simulates a waterfall-like pipeline with four stages (i.e., basic information synchronization, code review, code alignment, and document) and set up a code review team with six agents of different characters (i.e., user, CEO, CPO, CTO, coder, and reviewer). In the basic information synchronization phase, CEO, CPO, and coder agents analyze the input modality and programming language. After that, the coder and reviewer agents collaborate to conduct code review and produce the analysis report. In the code alignment phase, the coder and reviewer agents continue to revise the code based on the analysis reports. Finally, in the document phase, the CEO,

TABLE 5: Existing LLM-based Agents for Code Review

Agents	Multi-Agent Roles	Review Target			
		Consistency	Vulnerability	Code Smell	Code Optimization
CodeAgent [119]	User, CEO, CPO, CTO, Coder, Reviewer	✓	✓	✓	✓
Rasheed <i>et al.</i> [120]	Code Review, Bug Report, Code Smell, Code Optimization Agent		✓	✓	✓
ICAA [105]	Context & Prompt Incubation Agent, Consistency Checking Agent, Report Agent	✓			
CORE [121]	Proposer LLM, Ranker LLM				✓

CTO, and coder agents cooperate to document the holistic code review process. Experimental results demonstrate the effectiveness and efficiency of CodeAgent in various code review tasks, including consistency analysis, vulnerability analysis, format analysis, and code revision.

Different from CodeAgent which constructs a team that can conduct various code review tasks, Rasheed *et al.* [120] design an approach with each agent specialized for a single code review task individually. Notably, it proposes four agents including the code review agent, bug report agent, code smell agent, and code optimization agent. Each agent is trained on relevant Github data and evaluated on 10 AI-based projects. The results demonstrate the potential of applying multi-agent systems in the code review task.

ICAA [105] designs a multi-agent system to identify code-intention inconsistencies. It first uses the Context & Prompt Incubation Agent to collect necessary information from the code repository through a thinking-decision-action loop. The Consistency Checking Agent will then analyze collected information and identify inconsistencies, which will be handed over to the Report Agent to form a final report.

CORE [121] designs a system with two agents along with traditional static analysis tools to fix code quality issues automatically. Specifically, the Proposer agent takes the static analysis report, the suspicious file, and the issue documentation from language-specific static analysis tools (*e.g.*, CodeQL [114]) and the tool provider (*e.g.*, the QA team, and proposes candidate revisions for each suspicious file. After that, static analysis tools will prune revisions that still have issues, while the rest will be scored and re-ranked based on their likelihood of acceptance by the Ranker agent.

4.4 Testing

Software testing is essential for software quality assurance. LLMs have demonstrated promising proficiency in test generation, including generating test code, test inputs, and test oracles. However, generating high-quality tests in practice can be challenging, as the generated tests should not only be syntactically and semantically correct (*i.e.*, both the inputs and oracles should satisfy the specification of the software under test) but also be sufficient (*i.e.*, the tests should cover as many states of the software under test as possible). As shown by previous work [122], the tests generated by standalone LLMs still exhibit correctness issues (*i.e.*, compilation errors, run-time errors, and oracle issues) and unsatisfactory coverage. Therefore, researchers build LLM-based agents to extend the capabilities of standalone LLMs in test genera-

tion. We organize these works based on the test levels (*i.e.*, unit testing and system testing).

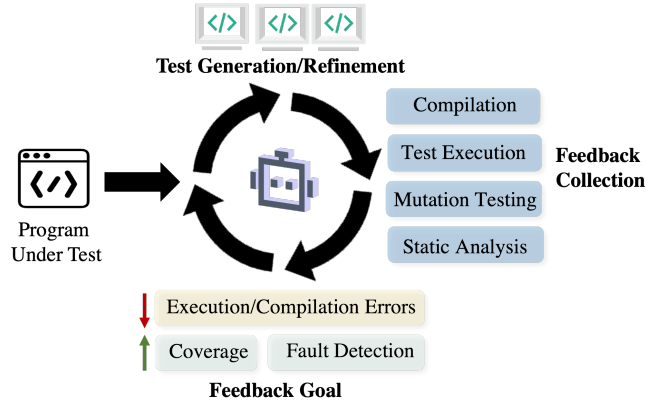


Fig. 9: Pipeline of LLM-based Agents for Unit Testing

4.4.1 Unit Testing

Unit testing checks the isolated and small unit (*e.g.*, method or class) in the software under test, which helps quickly identify and localize the bugs, especially for complicated software systems. Yuan *et al.* [122] perform a study showing the potentials of LLMs (*e.g.*, ChatGPT) in generating unit tests with decent readability and usability. However, the unit tests generated by standalone LLMs still exhibit compilation/execution errors and limited coverage. Therefore, recent works have built LLM-based agents that primarily extend standalone LLMs by iteratively refining the generated unit tests towards better correctness, coverage, and fault detection capabilities. Table 6 summarizes the existing LLM-based agents for unit test generation, and Figure 9 illustrates their common pipeline.

A. Iterative Refinement to Fix Compilation/Execution Errors. The test cases directly generated by LLMs can exhibit compilation or execution errors. Therefore, inspired by program repair [128], LLM-based agents further eliminate such errors by iteratively collecting the error messages and fixing the buggy test code [122]–[124]. For example, TestPilot [123] generates tests by constructing detailed prompts, including the function signature, implementation, documentation, and usage examples; it reflects on the feedback of failing tests and error messages to refine prompts and generate corrective tests iteratively. ChatTester [122] leverages LLMs to understand the intention of focal methods, and then generates a corresponding unit test; in addition, the iterative

TABLE 6: Existing LLM-based Agents for Unit Testing

Agents	Multi-Agent	Feedback Goal	Feedback Source	Target Language
ChatTester [122]	×	Reduce compilation/execution errors	Error messages	Java, Python
TestPilot [123]	×	Reduce compilation/execution errors	Error messages	JavaScript
ChatUniTest [124]	×	Reduce compilation/execution errors	Error messages	Java
TELPA [125]	×	Increase coverage	Program analysis results	Python
CoverUp [126]	×	Increase coverage	Execution results & Coverage	Python
MuTAP [127]	×	Enhance fault detection	Surviving mutants	Python

test refiner of ChatTester performs a more fine-grained refinement than TestPilot, which analyzes error messages and leverages static analysis tools to localize the buggy code for the next iteration refinement. Similarly, ChatUniTest [124] adopts a generation-validation-repair mechanism to refine the tests.

B. Iterative Refinement to Increase Coverage. CoverUp [126] is an LLM-powered test generation system designed to achieve high coverage rates. It dissects the task by segmenting the source code and employs the SlipCover tool [129] to conduct a detailed coverage analysis. Using an iterative approach, CoverUp refines its prompts to focus on areas of the code that lack coverage, thereby enhancing the overall quality and comprehensiveness of the generated test suite. Yang *et al.* [125] propose TELPA to enhance the coverage of hard-to-reach branches in software testing. Its methodology involves a two-pronged program analysis: backward and forward method invocation analysis, for a better understanding of the methods with uncovered branches. TELPA also employs counter-example sampling to guide LLMs toward generating novel tests that diverge from ineffective ones. The feedback-based process refines these tests iteratively through a CoT strategy, further improving coverage.

C. Iterative Refinement to Increase Fault Detection Capabilities. MuTAP [127] is a single LLM-based agent system that aims at generating unit tests of better bug detection capabilities with the feedback of mutation testing. It employs prompt augmentation with surviving mutants and refining steps to correct syntax and intended behavior. During each iteration, the LLM first generates initial test cases and self-refines their syntax errors and wrong behaviors, with the help of the Python parsing tool executing the tests. Then the tests run against the mutated programs, while the surviving mutants serve as feedback to direct the LLM in improving the test cases.

4.4.2 System Testing

System testing is a comprehensive process that assesses an integrated software system/component to guarantee that it fulfills its specification and operates as intended across diverse settings. For example, fuzzing testing and GUI (Graphical User Interface) testing are common testing paradigms at the system level. Leveraging LLMs for system testing can be challenging, as generating valid and effective system-level test cases should satisfy the constraints that are contained implicitly and explicitly in the specifications or domain knowledge of the software system under test. LLM-based agents are designed to better incorporate the domain knowledge of the software system under test compared to

generating system-level tests via standalone LLMs. We then organize these works according to the software systems under test. Table 7 summarizes the existing agents for different software systems.

A. OS Kernel. KernelGPT [130] is an LLM-based agent for kernel fuzzing. The analysis agent serves as the brain to automatically generate driver syscall specifications. It identifies drivers first and iteratively completes each specification component, during which it can determine whether additional information is needed based on its previous memory. A code extractor (implemented using LLVM toolchain [132]) is invoked to parse the kernel codebase for the provision of source code information. KernelGPT finally invokes the Syzkaller tool [131] (*e.g.*, `syz-extract`) and receives its feedback messages to validate and correct the generated syscall specifications iteratively.

B. Compiler. WhiteFox [133] encompasses two LLM-based agents working together, an analysis agent and a generation agent. The former examines the low-level optimization source code and produces requirements on the high-level test programs that can trigger the optimizations, while the latter crafts test programs based on summarized requirements. The generation agent further incorporates tests that have successfully triggered optimizations as feedback during the iterative process, thereby producing more satisfactory tests. LLM4CBI [134] is a single agent that aims at isolating compiler bugs by generating test cases of better fault detection capabilities. The agent utilizes tools to collect static information about the program (*e.g.*, `srcSlice` [136] for data flow) to construct precise prompts to guide the LLM for program mutation. The memorized component records meaningful prompts and selects better ones to instruct LLMs to generate variants. The generated programs undergo validation by static analysis (*e.g.*, the `Frama-C` tool [138]), and the feedback helps LLMs to avoid the same mistakes. The final test cases are used to identify suspicious files with spectrum-based fault localization techniques.

C. Mobile Applications. LLM-based agents are proposed to automate the testing process of mobile applications, including GUI testing, bug replay, and user acceptance testing.

C.1: GUI Testing. Some agents are developed to execute GUI testing for mobile applications. GUI testing is a commonly used software testing method aimed at verifying whether the user interface meets service specifications and user requirements. Previous LLM-based GUI testing approaches lack adequate autonomy, long-term planning, and coherence [154], [155]. The emergence of LLM-based agents enables GUI testing to focus more on higher-level test objectives [139], [144], [145], such as clear task objectives, without

TABLE 7: Existing LLM-based Agents for System Testing

Software System	Agents	Multi-Agent	Tool		Output
			Tool Category	Specific Tools	
OS Kernel	KernelGPT [130]	×	Static Analysis	syz-extract [131] LLVM Toolchain [132]	Syzkaller Specifications
Compiler	WhiteFox [133]	✓	-	-	Test Cases
	LLM4CBI [134]	×	Static Analysis	OClint [135] srcSlice [136] Gcov [137] Frama-C [138]	Mutated Programs
Mobile App	GPTDroid [139]	×	Execution Environment	VirtualBox [140] pyvbox [141] Android UIAutomator [142] Android Debug Bridge [143]	Test Scripts
	DroidAgent [144]	✓	Custom Toolkit	Navigation Action Toolkit	Test Scripts
	AXNav [145]	✓	Custom Toolkit	Navigation Action Toolkit	Bug Replay Video
	AdbGPT [15]	×	Execution Environment	Genymotion [146] Android UIAutomator2 [147] Android Debug Bridge [143]	Bug Replay Steps
	XUAT-Copilot [148]	✓	-	-	Test Scripts
Web App	RESTSpecIT [149]	×	-	-	HTTP Requests
Universal	Fuzz4All [150]	✓	-	-	Test Cases
	PentestGPT [151]	✓	Testing Tool	Metasploit [152]	Test Operations
	Fang <i>et al.</i> [153]	×	Custom Toolkit	Web Browsing Tool File Creation and Editing Tool	Exploit Actions
		Execution Environment	Terminal Code Interpreter		

relying on specific GUI states. Liu *et al.* [139] propose a framework called GPTDroid, where the LLM iterates the entire process by perceiving GUI page information, generating test scripts in the form of Q&A, executing these scripts through tools, and receiving feedback from the application. GPTDroid keeps a long-term memory to retain testing knowledge, which would help to improve the reasoning process. The DroidAgent [144] framework employs multiple LLM-based agents coordinating through different memory modules and can set its own tasks according to the functionalities of the apps under test. It is composed of four LLM-based agents: planner, actor, observer, and reflector, each with specific roles and supported by memory modules that enable long-term planning and interaction with external tools. AXNav [145] is another multi-agent system designed for replaying accessibility tests on mobile apps. It includes the planner agent, the action agent, and the evaluation agent, which together form the LLM-based UI navigation system. These agents translate test instructions into executable steps, conduct tests on a cloud-based iOS device, and summarize the test results in a chaptered video annotated with potential issues in the application, respectively.

C.2: Bug Replay. For automating Android bug replay, Feng *et al.* [15] introduce AdbGPT. Equipped with the knowledge of Step-to-Reproduce (S2R) entity specifications (*i.e.*, predefined actions and action primitives), AdbGPT analyzes bug reports to translate identified entities into a sequence of actions for bug reproduction using the CoT strategy. It then perceives GUI states dynamically and maps the S2R entities to actual GUI events to replicate the reported

bug.

C.3: User Acceptance Testing. To increase the automation of the user acceptance testing process, Wang *et al.* [148] propose XUAT-Copilot. The system is primarily comprised of three LLM-based agents responsible for action planning, state checking, and parameter selection, as well as two additional modules for state awareness and case rewriting. These agents interact with the testing equipment collaboratively, making human-like decisions and generating action commands.

D. Web Applications. RESTful APIs are popular among web applications as they provide a standardized, stateless, and easily integrable means of communication that enhances scalability and performance through a resource-oriented approach. RESTSpecIT [149] leverages LLMs to automatically infer RESTful API specifications and conduct black-box testing. Given an API name, RESTSpecIT generates and mutates HTTP requests through a reflection loop. By sending these requests to the API endpoint, it analyzes the HTTP responses for inference and testing. The LLM uses valid requests as feedback to refine the mutations in each iteration. Requests are validated based on the status code and message of the returned response.

E. Universal Software Categories. Some agent systems are not designed with a task-specific workflow, enabling them to be universally applicable across various target software systems. Xia *et al.* [150] present Fuzz4All, the first universal LLM-based fuzzer for general and targeted fuzzing across multiple programming languages. For a higher cost-effectiveness ratio, Fuzz4All consists of two agents, (i) the distillation LLM for user input distillation and initial

TABLE 8: Existing LLM-based Agents for Fault Localization

Agents	Multi-Agent	Tools		Input Context	FL Granularity	Target Language
		Tool Category	Specific Tools			
AgentFL [160]	✓	Static Analysis	Tree-sitter [161]	Project Level	Method	Java
RCAgent [162]	✓	Custom Toolkit	Code Analysis Tool Log Analysis Tool Memory Retrieval Tool Information Collection Tools	Project Level	Component	Java, Python
AUTOFL [163]	×	Custom Toolkit	Repository Retrieval Tools	Project Level	Method	Java

prompt generation, and (ii) the generation LLM for fuzzing input generation. They are powered by LLMs with different capabilities. In the fuzzing loop, the generation LLM refers to the previously generated samples and dynamically adjusts its strategy, thereby producing diverse fuzzing inputs. Deng *et al.* [151] design a modular framework, PentestGPT, to conduct Penetration Testing. The system includes inference, generation, and parsing modules. With the planning strategy of Pentesting Task Tree (which is based on the cybersecurity attack tree [156]) and CoT methods, PentestGPT solves the problems of context loss and inaccurate instruction generation that may be encountered during automated penetration testing. Fang *et al.* [153] develop a benchmark consisting of 15 one-day vulnerabilities to assess the efficacy of their agent framework in exploiting such weaknesses, utilizing various LLM backbones. Their agents are imbued with an understanding of the Common Vulnerabilities and Exposures (CVE) descriptions and are capable of harnessing a suite of tools to facilitate the exploitation process. These tools include web browsing capabilities for navigation, web search functionalities for traversing web pages, as well as terminal and code interpreter access for the generation and execution of scripts.

4.5 Debugging

Software debugging typically includes two phases: *fault localization* [157] and *program repair* [158]. In particular, fault localization techniques aim at identifying buggy elements (*e.g.*, buggy statements or methods) of the program based on the buggy symptoms (*e.g.*, test failure information); then, based on the buggy elements identified in the fault localization phase, program repair techniques generate patches to fix the buggy code. In addition, recent works also propose *unified debugging* to bridge fault localization and program repair in a bidirectional way [159]. We then organize the works in LLM-based agents for debugging into three parts, *i.e.*, fault localization, program repair, and unified debugging.

4.5.1 Fault Localization

Learning-based fault localization has been widely studied before the era of LLMs, which typically trains deep learning models to predict the probability of each code element being buggy or not [164]. However, precisely identifying the buggy element in the software is challenging, given the large scale of the software systems as well as the massive and diverse error messages, which are often beyond the capabilities of standalone learning models including LLMs. Therefore, recent works build LLM-based agents, which incorporate multi-agents and tool usage to help LLMs tackle these challenges. Table 8 summarizes the existing LLM-based agents for fault localization, and Figure 10 illustrates their common pipeline.

A. Multi-agent Synergy. AgentFL [160] is a multi-agent system for project-level fault localization. The main insight of AgentFL is to scale up LLM-based fault localization to project-level code context via the synergy of multiple agents. The system consists of four distinct LLM-driven agents: test code reviewer, source code reviewer, software architect, and software test engineer. Each agent is customized with specialized tools and a unique set of expertise. With the four agents, AgentFL streamlines the project-level fault localization process by breaking it down into three phases: fault comprehension, codebase navigation, and fault confirmation.

RCAgent [162] is a multi-agent system for root cause analysis in industrial cloud settings. RCAgent includes two components: the controller agent and expert agents. The controller agent oversees the comprehensive thought-action-observation cycle, while the expert agents act for specialized tasks and can be utilized by the controller agent. A key-value store is employed for the controller agent to memorize the observation information (*e.g.*, logs and table entries) to help the decision-making, as well as handling the context length constraint. The expert agents perform code analysis and log analysis tasks respectively, and their summarized results are fed back to the controller agent as observation. The agents can invoke tools for information collection (*e.g.*, log data and repositories retrieval) or memory retrieval, through the function calling. Besides, a self-consistency mechanism is built to enhance the performance.

B. Tool Invocation. AUTOFL [163] is a single-agent system, which enhances standalone LLMs with tool invocation (*i.e.*, four specialized function calls) to better explore the repository. It first performs root cause explanation, invoking tools to oversee the source code repository for pertinent information, requiring only a single failing test and its failure stack. During this stage, it autonomously decides

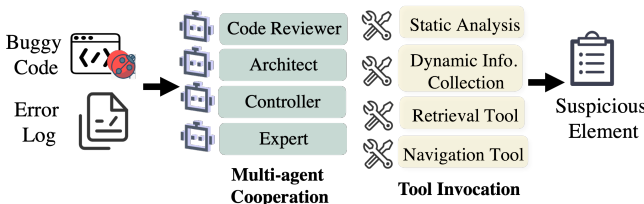


Fig. 10: Pipeline of LLM-based Agents for Fault Localization

TABLE 9: Existing LLM-based Agents for Program Repair

Agents	Multi-Agent	Feedback Source	Target Software	Benchmark	Fix Rate
ChatRepair [166]	×	Execution/Compilation	Java	Sampled Defects4J [167] & QuixBugs [168]	162/337 (Defects4J) 80/80 (QuixBugs)
CigaR [169]	×	Execution/Compilation	Java	Sampled Defects4J & HumanEval-Java [170]	69/267 (Defects4J) 102/162 (HumanEval)
RepairAgent [171]	×	Execution/Compilation	Java	Defects4J	164/835
AutoSD [172]	✓	Execution	Java/Python	Defects4J & BugsInPy [173] & Almost-Right HumanEval	189/835 (Defects4J) 187/200 (HumanEval)
ACFIX [174]	✓	Static Checking Model Debate	Smart Contract	Self-curated Dataset	112/118
Conversational APR [128]	×	Execution/Compilation	Java/Python	Sampled QuixBugs	59/60
FlakyDoctor [175]	✓	Execution Static Checking	Java	Sampled IDoFT [176] & DexFix dataset [177] & Sampled ODRRepair dataset [178]	245/419 (Implementation-Dependent Flakiness) 185/247 (Order-Dependent Flakiness)

whether to continue function calling or to terminate with the production of root cause explanation. Subsequently, a post-processing step is used to correlate the outputs with exact code elements, aiming at bug localization. In addition, AgentFL [160] and RCAGENT [162] also incorporate tool invocation (e.g., static analysis, dynamic instrument, and code base navigation) into their framework.

4.5.2 Program Repair

Fine-tuning and fixed prompting have been the most widely adopted paradigms for program repair techniques based on standalone LLMs. In particular, program repair is formulated as a translation problem [165] (i.e., translating the buggy code to correct code) or a generation problem [12] (e.g., infilling the correct code in the buggy code context). However, patches generated by LLMs in a single iteration are not always correct; they may fail to pass all the tests or may be overfitting to the test cases. Therefore, existing LLM-based agents all follow an iterative paradigm to refine patch generation based on the tool or model feedback in each iteration. Table 9 summarizes the existing LLM-based agents for program repair, and Figure 11 illustrates their common pipeline.

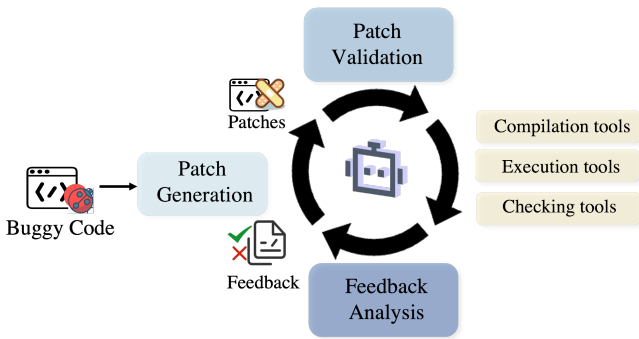


Fig. 11: Pipeline of LLM-based Agents for Program Repair

ChatRepair [128], [166] is the first automated approach to iteratively refine patch or program generation based on environmental feedback. In each iteration, ChatRepair leverages tools to compile and execute the generated patches, and generates new patches based on the compilation/execution feedback and also earlier patch attempts in the same session.

CigaR [169] is a similar agent system for function-level program repair. CigaR leverages feedback to refine its outputs iteratively and decides to reboot the repair process

when needed, ensuring the production of both plausible and diverse patches.

RepairAgent [171] adopts a more agentic design and improves the iterative refinement procedure by making the hard-coded feedback mechanism more flexible. In particular, RepairAgent allows the LLM itself to decide when or which tool to invoke. In addition to the basic tools (e.g., compilation tools and test execution tools) used in previous work, RepairAgent further includes tools for code reading, codebase searching, and hypotheses stating/discarding.

AutoSD [172] is a multi-agent system that iteratively fixes the buggy program via simulating the scientific debugging [179]. AutoSD includes four components: LLM-based hypothesis generator, execution-based validator, LLM-based conclusion maker, and LLM-based fixer. In each iteration, the generator first generates a hypothesis about the bug, then invokes the debugger tool for hypothesis validation; the conclusion maker further identifies whether the hypothesis is rejected or not; and the fixer finally returns potential patches with explanations.

ACFIX [174] is a multi-agent system for fixing the access control vulnerabilities in smart contracts. By specializing LLMs with different roles, ACFIX includes a Role-based Access Control (RBAC) mechanism identifier, a role-permission pair identifier, a patch generator, and a validator. In particular, the validator checks the validity of the generated patches with both tool feedback (static grammar rule checking) and model feedback (multi-agent debate process). The feedback is further provided to iteratively refine the patch.

FlakyDoctor [175] is an agent to repair flaky tests. It takes into consideration test execution results and the location of test failures. Following this, the system generates targeted repairs and tests them for validation. This process is iterative, with the aim of continually refining the repairs until the issue of test flakiness is resolved.

4.5.3 Unified Debugging

Instead of tackling fault localization or program repair as isolated phases, unified debugging techniques treat them as a unified procedure, which leverages the outputs of each phase to refine each other. In particular, traditional unified debugging techniques [180]–[182] primarily pre-define heuristic rules to refine fault localization based on the patch validation results during program repair. Recently, LLM-based agents have enhanced traditional unified debugging techniques with more flexibility by leveraging LLMs to

TABLE 10: Existing LLM-based Agents for End-to-end Software Development

Agents	Multi-Agent	Process Model	Roles Creation	Collaboration Mode	Communication Protocol
Self-Collaboration [4]	✓	Waterfall	Pre-defined	Ordered	Natural Language
Low-code LLM [183]	✓	-	Pre-defined	Ordered	Natural Language
Prompt Sapper [184]	×	-	Pre-defined	Ordered	Natural Language
Talebirad <i>et al.</i> [185]	✓	-	Dynamic	Ordered	Natural Language
ChatDev [186]	✓	Waterfall	Pre-defined	Dual-role	Natural Language
MetaGPT [187]	✓	Waterfall	Pre-defined	Ordered	Structured
AgentVerse [188]	✓	-	Dynamic	Ordered	Natural Language
AutoAgents [189]	✓	-	Dynamic	Ordered	Natural Language
Qian <i>et al.</i> [190]	✓	-	Pre-defined	Ordered	Natural Language
AISD [191]	✓	Waterfall	Pre-defined	Ordered	Natural Language
LLM4PLC [192]	×	-	Pre-defined	Ordered	Natural Language
CodePori [193]	✓	-	Pre-defined	Ordered	Natural Language
LCG _{Waterfall} [194]	✓	Waterfall	Pre-defined	Ordered	Natural Language
LCG _{Scrum} [194]	×	Agile	Pre-defined	Debate	Natural Language
CodeS [195]	✓	-	Pre-defined	Ordered	Natural Language
Qian <i>et al.</i> [196]	✓	-	Pre-defined	Ordered	Natural Language
CTC [197]	✓	Waterfall	Pre-defined	Dual-role	Natural Language
AgileCoder [198]	✓	Agile	Pre-defined	Dual-role	Natural Language

comprehend, utilize, and unify the outputs of both fault localization and program repair.

FixAgent [159], a multi-agent system for unified debugging, enables end-to-end fault localization, bug repair, and bug analysis. Based on manual debugging (*e.g.*, rubber duck debugging), FixAgent uses agent specialization and synergy (*i.e.*, LLM localizer, LLM repairer, LLM crafter, and LLM revisitor) to incorporate key variable tracking and program context comprehension. FixAgent can fix 79 out of 80 bugs in the QuixBugs [168] benchmark.

LDB [75] is an agent for end-to-end fault localization and program repair. LDB divides the buggy program into basic blocks according to a control-flow graph, and leverages LLMs to detect the incorrect blocks with run-time execution values and then to provide refinement suggestions.

4.6 End-to-end Software Development

Given the high autonomy and the flexibility from multi-agent synergy, LLM-based agent systems can further tackle the end-to-end procedure of software development (*e.g.*, developing a Snake Game application from scratch) beyond an individual phase of software development. In particular, alike the real-world software development team, these agent systems can cover the entire software development life cycle (*i.e.*, requirements engineering, architecture design, code generation, and software quality assurance) by incorporating the synergy between multiple agents that are specialized with different roles and relevant expertise. Table 10 summarizes the existing LLM-based agents for end-to-end software development.

4.6.1 Software Development Process Model

Real-world software teams often follow classic software process models (*e.g.*, waterfall [199], incremental model [200], unified process model [201], and agile development [202]) to facilitate a more standardized software development life cycle. Existing LLM-based agents for end-to-end software

development also design their workflows according to common software process models, *e.g.*, waterfall process model and agile development. Figure 12 illustrates how existing LLM-based agents go through different process models.

A. Waterfall Process Model. The majority of existing LLM-based agent systems (*e.g.*, AISD [191], LCG [194], ChatDev [186], CTC [197], and Self-Collaboration [4]) follow the classic waterfall process model for software development. The traditional waterfall process model [199] is a linear and sequential software development workflow that divides the project into distinct phases, *i.e.*, requirements engineering, design, code implementation, testing, deployment, and maintenance. Once a phase is finished, the project moves forward to the next phase without iteration. Based on this process, some end-to-end software development agents [4], [187], [191], [194] further extend the traditional waterfall process by including iterations in specific phases to ensure the high quality of the generated content. For example, the results of the testing phase might be fed back to the developer agent to revise the generated code; MetaGPT [187] further integrates the waterfall model with human-like Standardized Operating Procedures (SOPs), which assign responsibilities to each role and standardize the intermediate outputs, promoting collaboration among different team members.

B. Agile Development. Some works explore the potential of LLM-based agents with the agile development, including Test-Driven-Development (TDD) [194] and Scrum [194], [198]. TDD prioritizes writing tests before the actual coding and fosters a cycle of writing test suites, implementing the code to pass the test suites, and concluding with a reflective phase to refinement. Scrum is an agile software development process model that breaks down software development into several sprints, achieving complex software systems through iterative updates. Experiments on function-level code generation benchmarks show that the Scrum model can achieve the best and most stable

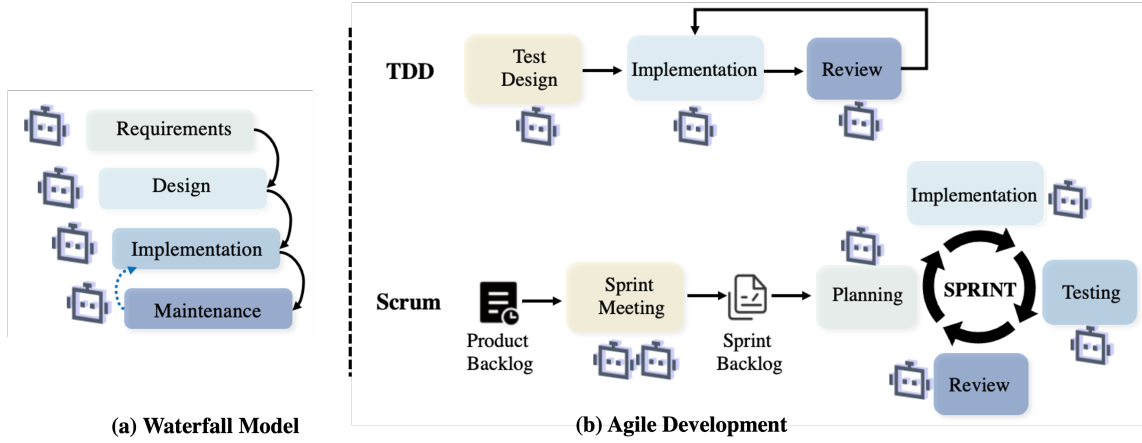


Fig. 12: Process Models Adopted by LLM-based Agents for End-to-end Software Development

performance, followed by the TDD model [194].

4.6.2 Role Specialization of Software Development Team

Imitating real-world software development teams, multi-agent systems for end-to-end software development often assign different roles to tackle specialized sub-tasks and collaborate throughout the software development life cycle.

A. Role Categories. The roles in existing agents are primarily designed by simulating real-world software development teams or specialized by the workflow.

A.1: Simulating Real-world Software Teams. Most end-to-end frameworks simulate the real-world software development teams and assign basic roles including managers (e.g., project managers or product managers), requirement analyzers, designers, developers, and quality assurance experts (e.g., software testers or code reviewers) to cover the entire pipeline of software development [4], [187], [191], [193], [194], [198]. In addition to these common roles, there are some special roles that can be assigned to tackle fine-grained tasks. For example, a Scrum master role is also included into the requirements analysis and planning tasks for the agents with the Scrum workflow [194], [198]; and CEO/CTO roles are also included in some agents to complete the design task [186], [197]; in addition, there are special supervisor roles in some agents to ensure the smooth progress of the collaboration (e.g., providing coordination or critiques), such as the oracle roles in previous work [185] and the action observer in AutoAgents [189]. The detailed categories of roles in existing agents are discussed in Section 5.2.1

A.2: Specialized by Workflow. Instead of simulating the real-world development teams, some agents break down roles according to the agent framework workflow. For example, CodeS [195] decomposes the complex code generation task into the implementation of repository, file, and method layers, and sets up the roles of RepoSketcher, FileSketcher, and SketchFiller. Co-Learning [190] and its subsequent work [196] abstract the code generation process into instruction-response pairs, thus only setting up the roles of instructor and assistant.

B. Role Creation. The roles in multi-agent systems are either created in a pre-defined way or in a dynamic way.

B.1: Pre-defined. The majority of specialized roles are pre-defined by the agent framework [4], [183], [186], [187], [190],

[191], [193]–[198]. In other words, the roles are fixed by the agent for each task.

B.2: Dynamic Creation. In addition to the pre-defined roles, some agents assign roles in a dynamic way, which can equip multi-agent systems with more flexibility. For example, AutoAgents [189] designs a drafting stage that aims at determining the roles of the multi-agent group via the communication between two meta agents: the planner and the agent observer; in addition, AgentVerse [188] sets up a group of different roles through an expert recruitment stage. Talebirad *et al.* [185] propose a novel framework and enable an agent to spawn additional agents to the system. Such dynamic strategies aim at creating roles in a more diverse and flexible way.

4.6.3 Collaboration Mechanism in Multi-agent

Within the multi-agent systems for end-to-end software development, it is essential to schedule how each agent is coordinating with each other. We then discuss the collaboration mode and the communication protocol adopted in existing agents.

A. Collaboration Mode. In particular, there are mainly two collaboration modes in multi-agent systems for end-to-end software development, *i.e.*, the ordered mode and the unordered mode.

A.1: Ordered Mode. It is a sequential mode wherein each agent makes decisions independently and uses optional feedback mechanisms to improve quality of their generated content. Ordered mode is the most prevalent collaboration mode adopted in existing agent systems [4], [186]–[193], [195]–[198]. Previous research [188] suggests that this approach is more suitable for software development, as it focuses on producing only the final refined decision.

A.2: Unordered Mode. It primarily introduces an unordered debate mechanism, where multiple agents present their opinions separately and equally, with the ultimate decision reached via summarization. For example, LCG [194] involves the sprint meeting in the Scrum process model, wherein participating agents propose their opinions and share them with all other agents through a buffer. The Scrum master ultimately summarizes and extracts user stories. Besides, some works adopt a multi-role mechanism

TABLE 11: Benchmarks for End-to-end Software Development

Benchmarks	#Tasks	Input Scale	Output Scale	Language	Evaluated Agents
SRDD [186]	1,200	Software Description (55 words)	Multiple Files	Python	[186], [190], [196], [197]
CAASD [191]	72	Software Description (50 words)	Multiple Files	Python	[186], [187], [191]
SoftwareDev [187]	70	Software Description (30 words)	Multiple Files	Python	[186]–[188]
SketchEval [195]	19	README (421 words)	Structured Multiple Files	Python	[186], [195]
ProjectDev [198]	14	Software Description (262 words)	Multiple Files	Python	[186], [187], [198]
HumanEval [170] HumanEval-ET [203]	164	Function Description (68 words)	Single Function	Python	[4], [187]–[189], [193], [194], [198]
MBPP [204] MBPP-ET [203]	974	Function Description (15 words)	Single Function	Python	[4], [187], [193], [194], [198]

TABLE 12: Metrics Used in Evaluating Agents for End-to-end Software Development

Category	Metrics	Used Agents
Execution Validation	Pass Rate PassK Executability #Errors	[4], [186], [187] [188]–[190] [193], [194], [196] [197], [198]
Similarity	SketchBLEU [195] Cosine Distance	[186], [190], [195] [196], [197]
Costs	Running Time Token Usage Expenses #Sprints	[187], [198]
Manual Efforts	Human Revision Costs	[187]
Generated Code Scale	Line of Code Code Files Completeness	[186], [187], [190] [196], [197]

to ensure the generation quality [186], [190], [196]–[198]. Instead of passing intermediate results from one agent to another, these works assign two agents to work together to complete a sub-task in a communication way. For example, in ChatDev [186] and CTC [197], agents with different roles will collaborate to complete specific tasks (e.g., the system design is achieved through communication between the CEO and CTO).

B. Communication Protocol. Within the multi-agent systems, agents communicate with other agents to exchange information. In particular, there are two communication protocols, *i.e.*, the pure natural language and the structured communication.

B.1: Natural Language. The most common communication protocol is direct dialogue [4], [186], [188], [189], [191], [193], [194], which leverages natural language to exchange information. This approach allows for flexible expression of intent and is close to human communication.

B.2: Structured. Some agents (e.g., MetaGPT [187]) structure communication by having agents exchange documents and diagrams instead of relying solely on dialogue, as pure natural language may be insufficient for solving complex tasks due to distortion in multi-turn communication [187].

4.6.4 Agent Evaluation

Given the complexity of end-to-end software development, researchers further build diverse benchmarks and metrics for a comprehensive evaluation.

A. Benchmarks. Table 11 summarizes the benchmarks used for evaluating existing LLM-based agents for end-to-end software development. In particular, we can observe that there are still a large number (e.g., 7) of papers using the classic code generation benchmarks (e.g., HumanEval [170] or MBPP [204]) for evaluating end-to-end software development. Although these traditional code benchmarks can represent end-to-end software development to some extent, they still involve simplified, small-scale development tasks (*i.e.*, input of short function descriptions and output of a single function, as shown in Table 11). In addition, there are five more complicated benchmarks that aim at simulating the end-to-end software development, *i.e.*, SRDD [186], [190], [196], [197], CAASD [191], SoftwareDev [187], SketchEval [195], and ProjectDev [198]. The tasks in these benchmarks include more complicated and longer requirement description (e.g., the average length of software description in ProjectDev is 262 words), and their expected outputs are supposed to contain multiple files. In particular, the benchmark SketchEval is built upon the real-world GitHub repositories, and its input descriptions are extracted from the README file of the software while its output expects multiple files that are organized in a repository structure.

B. Metrics. Table 12 summarizes the metrics used for evaluating existing LLM-based agents for end-to-end software development. In fact, given the difficulty of generating complicated program, it can be possible that the generated program cannot perfectly pass the tests. Therefore, in addition to the common metrics (e.g., Pass Rate or Pass@K) that execute the generated program for validation, there are multiple dimensions for assessing how existing agents perform in end-to-end to software development. In particular, there are (i) the similarity metrics between the generated program and the ground truth (e.g., SketchBLEU [195] measures the structure similarity), (ii) the costs of executing or generating the program, (iii) the manual efforts to further refine the generated program, and (iv) the scale of the generated

TABLE 13: Existing LLM-based Agents for End-to-end Software Maintenance

Agents	Multi-Agent	Phases						
		Preprocessing	Issue Reprod.	Issue Localization	Task Decomp.	Patch Generation	Patch Verification	Ranking
MAGIS [205]	✓	×	×	Retrieval-based	×	w/ local context	Code Review	×
AUTOCODEROVER [206]	✓	×	×	Navigation/Spectrum-based	×	w/ cross-file context	Static Check	×
SWE-agent [207]	×	×	✓	Navigation-based	×	w/ local context	Static Check	×
CodeR [208]	✓	Plan Selection	✓	Spectrum-based	×	w/ cross-file context	Dynamic Check	×
RepoUnderstander [209]	✓	Knowledge Graph Const.	×	Simulation	✓	w/ cross-file context	Static Check	×
MASAI [210]	✓	Test Template Generation	✓	Navigation-based	✓	w/ local context	Static/Dynamic Check	✓
Agentless [211]	×	Repository Tree Const.	×	Navigation-based	×	w/ local context	Static/Dynamic Check	✓

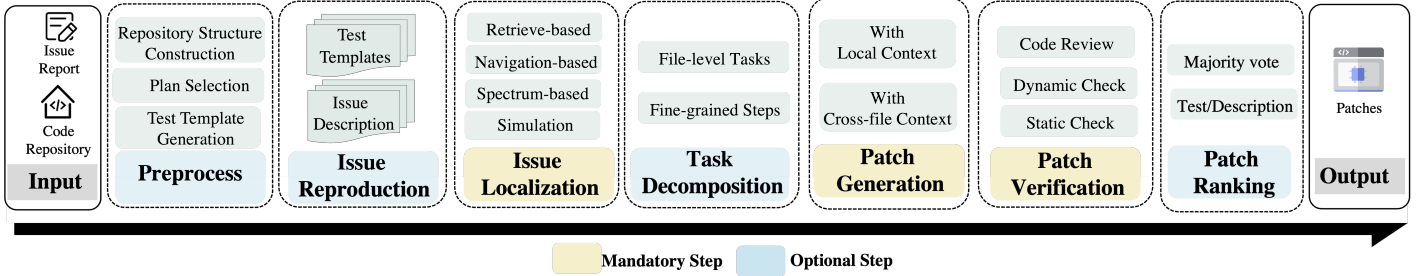


Fig. 13: Pipeline of LLM-based Agents for End-to-end Software Maintenance

program.

4.7 End-to-end Software Maintenance

Software systems undergo maintenance as requirements continuously change (*i.e.*, adding, deleting, or modifying features) or unexpected software behaviors arise. In practice, users report unsatisfactory behaviors that they encounter; developers then diagnose the reported issues and modify the software to fix them. Such an end-to-end software maintenance process can be time-consuming and labor-intensive in practice, as it involves multiple phases including understanding user-reported issues, localizing code for maintenance, and precisely editing code to address issues. Recently, there has been an increasing number of multi-agent systems aiming at automatically solving issues of real-world software projects. Table 13 summarizes the characteristics of these agents.

4.7.1 Common Pipeline

Figure 13 illustrates the pipeline of existing LLM-based agent systems for end-to-end software maintenance. Basically all of the existing agents follow a common pipeline of three phases, *i.e.*, *issue localization*, *patch generation*, and *patch verification*, where different agents incorporate different strategies to tackle each phase. In addition, some agents further include additional phases, *i.e.*, *preprocessing*, *issue reproduction*, *issue localization*, *task decomposition*, *patch generation*, *patch verification*, and *patch ranking*.

A. Preprocessing. To better understand the whole repository, some agents first perform preprocessing to prepare pre-knowledge before the entire procedure. The agent system RepoUnderstander [209] constructs a knowledge graph of the entire code repository to facilitate the subsequent process of issue localization. Meanwhile, Agentless [211], which is simplistic and less agentic, simply turns the whole project into a tree-like structure that demonstrates all directories and files of the repository in a hierarchical format, which

facilitates the issue localization phase. In CODER [208], a manager agent first chooses a plan from several workflows pre-defined by human experts. In MASAI [210], the test template generator is used to analyze the testing setup of the repository and generate a test template with the running command, which further serves as an example for the following issue reproduction phase.

B. Issue Reproduction. A test script that triggers the unexpected behaviors users encounter is essential for issue resolution. It not only helps with issue localization but also serves as the verification criterion for patch correctness. However, in practice, users do not always provide such reproduction tests when they report issues; and such reproduction tests are often added by developers after they fix the buggy software. Therefore, some agents design the issue reproduction phase that aims at generating the test script that can trigger the unexpected behaviors encountered by users. For example, SWE-agent [207] and CodeR [208] directly leverage LLMs to generate reproduction tests based on issue descriptions when there is no existing reproduction script in issue descriptions. However, generating reproduction tests can be challenging, as the tests must be executable and ideally should fail on the buggy software version while passing on the fixed version. Therefore, to increase the success rate of issue reproduction, the multi-agent system MASAI [210] includes a two-stage approach for issue reproduction, which first investigates the test framework and existing tests for generating a sample test template (generated in the preprocessing phase) and then uses the template as a demonstration to create the reproduction script.

C. Issue Localization. Issue localization is one of the most important phases where the agents are supposed to precisely identify the code elements (*e.g.*, classes, methods, or code blocks) that are related to issues and should be edited. We then summarize the common localization strategies used in existing LLM-based agents.

C.1: Retrieval-based Localization. All the existing agents

use the retrieval-based strategy as the basis for issue localization, which identifies the relevant code elements based on their similarity with issue descriptions. For example, in MAGIS [205], all code files are compared to issue descriptions via BM25 [212], and the Top-K relevant code files are selected as the potential issue locations.

However, only using retrieval-based strategy can be insufficient and often at coarse granularities (*e.g.*, files). Therefore, some agents further extend the basic retrieval with other strategies, *e.g.*, navigation-based localization, spectrum-based localization, and simulation-based localization.

C.2: Navigation-based Localization. This strategy provides agents with a set of code search actions that can navigate through the entire code repository to check all directories and files. For example, the issue localization phase of SWE-agent [207] uses several pre-defined search-based interfaces to locate the target directories or files and then to view the code snippets in the target file through scrolling interfaces. Similarly, MASAI [210] assigns an edit localizer that can navigate the repository to find the related code snippets. AUTOCODEROVER [206] designs a stratified context retrieval process, which allows the LLM itself to decide whether to further refine the location based on the current context, thus forming an iterative navigation process. Agentless [211] provides the hierarchical structure of the target repository and instructs the LLM to gradually localize files, classes, functions, and concrete edit locations.

C.3 Spectrum-based Localization. Some agents integrate the traditional fault localization approaches, especially spectrum-based fault localization techniques [157], which calculate the suspiciousness score of code elements based on their coverage of failed tests and passing tests. For example, AUTOCODEROVER [206] explores spectrum-based techniques for issue localization by using the ground-truth reproduction tests provided in SWE-bench Lite, which improves issue resolution rate from 17.00% to 20.33%. While AUTOCODEROVER explores the spectrum-based fault localization in an ideal case (as the ground-truth reproduction tests are not always available in practice), CODER [208] investigates the improvement of the spectrum-based fault localization in a more practical setting by using tests generated in the issue reproduction phase. In particular, CODER calculates the suspiciousness scores based on the coverage of the reproduction tests, and then combines them with the basic retrieval-based strategy (*i.e.*, BM25 similarity between code elements and issue descriptions) via weighted computation.

C.4 Simulation. Simulation is a special technique adopted by RepoUnderstander [209] for issue localization. It applies the classic Monte Carlo Tree Search algorithm. By recursively incorporating nodes of the high BM25 score with the issue, it evaluates and ranks the most relevant paths in the repository knowledge graph. The collected code is then summarized for issue localization.

D. Task Decomposition. Before generating patches, some agents decompose the task into more fine-grained sub-tasks. For instance, in MAGIS [205], its manager agent breaks down the issue into file-level tasks and delegates them to a newly-formed development team; similarly, in RepoUnderstander [209], its summary agent summarizes

the collected code and issue description, and then outlines the fine-grained steps for issue resolution.

E. Patch Generation. In this phase, the agents generate patches for the localized suspicious code elements. The input context of this phase typically includes the issue/task description and the suspicious code elements for modification [205], [207], [210], [211]. In addition, some agents (*e.g.*, AUTOCODEROVER [206], CodeR [208], and RepoUnderstander [209]) further refine the input contexts by including relevant cross-file code contexts that are collected by retrieval APIs.

F. Patch Verification. Agents further verify the correctness of the generated patches, which is challenging as the reproduction tests are not always available in practice. Therefore, agents incorporate different verification strategies.

F.1: Code Review. Some agents (*e.g.*, MAGIS [205]) design a quality assurance agent to review the quality of generated patches.

F.2: Static Checking. Some agents (*e.g.*, AUTOCODEROVER [206], RepoUnderstander [209], MASAI [210], Agentless [211], and SWE-agent [207]) use static checking approaches to assess the syntactic correctness, indentation, and compatibility of the generated patch with the repository environment.

F.3: Dynamic Checking. Since the static checking cannot find the semantic violation of the patches, some agents (*e.g.*, CodeR [208] and MASAI [210]) further perform dynamic checking by executing the reproduction test on the patch. The patch that passes the reproduction test can be considered as effectively resolving the issue. In particular, existing reproduction tests are reused (if available); otherwise, reproduction tests generated during the issue reproduction phase are used. Agentless [211] also implements a dynamic checking approach by conducting regression testing to filter out incorrect candidate patches.

G. Patch Ranking. Since the patch verification phase can sometimes be insufficient for filtering out all the incorrect patches, some agents further include a patch ranking phase to identify the patch with the highest probability of being correct. For example, in MASAI [210], a ranker agent is responsible for ranking all potential patches based on the issue description and reproduction tests; In Agentless [211], all patches are normalized and re-ranked based on the number of occurrences with the majority voting strategy.

4.7.2 Benchmarks

To evaluate how LLM-based agents tackle end-to-end software maintenance, researchers build benchmarks from real-world Github issues, including SWE-bench [213], SWE-bench Lite [214], SWE-bench Lite-S [211], and SWE-bench Verified [215]. Table 14 summarizes the evolution timeline of existing benchmarks for end-to-end software maintenance.

SWE-bench [213] is the first benchmark for end-to-end software maintenance, which consists of 2,294 real-world GitHub issues across 12 popular Python repositories. Each task in SWE-bench includes an original text from a GitHub issue (*i.e.*, the issue description or problem statement), the entire code repository, the execution environment (*i.e.*, Docker environment), and validation tests (*i.e.*, tests that are hidden from the evaluated agents).

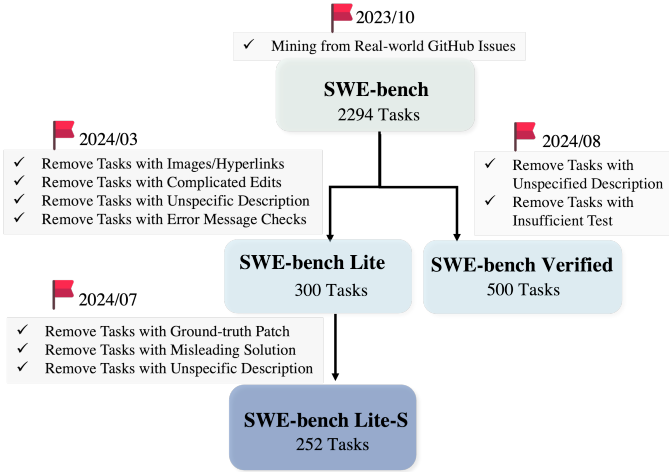


Fig. 14: Benchmark Evolution in Software Maintenance

However, the full SWE-bench benchmark can take too much evaluation costs and it contains particularly difficult or problematic tasks [211], which can underestimate the evaluation of LLM-based agents. Therefore, researchers have dedicated lots of manual efforts to identifying high-quality tasks with reasonable difficulty, self-contained information, informative issue descriptions, and sufficient evaluation tests. For example, SWE-bench Lite [214] is a subset of SWE-bench that manually removes tasks requiring complicated edits (*e.g.*, editing more than one file) or the tasks including images or hyperlinks; SWE-bench Lite-S [211] removes tasks that contain exact patches, misleading solutions, or insufficient information in the issue descriptions; similarly, SWE-bench Verified [215] removes cases with unspecified descriptions or insufficient tests.

5 ANALYSIS FROM AGENT PERSPECTIVE

This section organizes the collected papers from the perspective of agents. Specifically, Section 5.1 summarizes the components of existing LLM-based agents for SE; Section 5.2 focuses on existing multi-agent systems for SE by summarizing their roles and collaboration mechanisms; and Section 5.3 summarizes how humans coordinate with agents for SE.

5.1 Agent Framework

Based on the common framework of LLM-based agents [21], [26], [28], this section summarizes the common paradigms of the planning, perception, memory, and action components in existing LLM-based agents for SE.

5.1.1 Planning

In SE, intricate tasks such as development and maintenance activities necessitate the orchestrated efforts of various agents through multiple iterative cycles. Therefore, planning is an essential component for agent systems by meticulously delineating task sequences and strategically scheduling agents to ensure the seamless progression of the SE process. Figure 15 presents the taxonomy of the planning components in existing LLM-based agents for SE.

A. Single Planner vs. Multiple Planners. In LLM-based agent systems, planning is typically handled by a specialized agent [4], [61], [91], [98], [105], [144], [145], [183], [187], [191], [193] or as a core responsibility of an individual agent [76], [82], [85], [86], [148], [162], [192]. Some works use the function-calling interface [216] provided by GPT-3.5 [217] or GPT-4 [218], handing over the planning task to high-performance models [82]. However, given the pivotal role that planning plays in influencing subsequent action steps, some works incorporate a collaborative approach among several agents to further enhance the accuracy and practicality of the plans formulated [91], [186], [189], [195], [197], [198], [205].

B. Single-turn Planning vs. Multi-turn Planning. The fundamental planning strategy is to craft a holistic plan from the very beginning meticulously and then proceed to implement it in successive rounds [4], [61], [82], [86], [91], [98], [105], [183], [186], [187], [191]–[193], [195], [197], [198], [205]. Further, many SE agents have adopted a ReAct-like [219] architecture, which implements a multi-turn planning mechanism wherein the next-round actions will not be determined until receiving the environmental feedback from the previous round. This form allows for dynamic revision and expansion of the plan, enabling it to adapt to more flexible task scenarios, such as issue resolution [76], [210], iterative code generation [82], [85], mobile app testing [144], [145], [148], among others [162].

C. Single-path Planning vs. Multi-path Planning. Most LLM-based agents use single-path planning strategies, *i.e.*, they plan and execute tasks in a linear manner [4], [61], [86], [105], [151], [186], [187], [191]–[195], [198], [205], [220]. However, agents inherit the randomness from the backbone LLMs, leading to fluctuations in task decomposition. Some approaches improve upon single-path planning by using feedback from each round to dynamically plan the next round of actions [82], [85], [88], [97], [144], [145], [148], [162], [210]. Although these dynamic strategies are still single-path, they offer considerable flexibility due to their ability to be adjusted based on progress and execution outcomes. Another approach to address this issue is to design a multi-path planning strategy, which instructs the agents to generate or simulate multiple plans, and select [76], switch [98], or aggregate [197] the optimal paths for execution.

D. Plan Representation. The plan can be exhibited in different forms, including natural language descriptions, semi-structured representations, or graphs.

- *Natural Language.* Most agents describe the plan in natural language, especially as a list of procedural steps [4], [86], [91], [98], [105], [183] or features to be implemented [187], [191], [194], [198].
- *Semi-structured.* The agent system AXNav [145] represents the action list in JSON format; and some code-generating agent systems directly output the code skeleton [61], [97], [192], [195] or present the plan as executable code [205], which can be seen as a special plan form in SE tasks.
- *Graph.* Some agents model the plan as a graph to facilitate the expansion and traceability of execution paths [76], [88], [151].

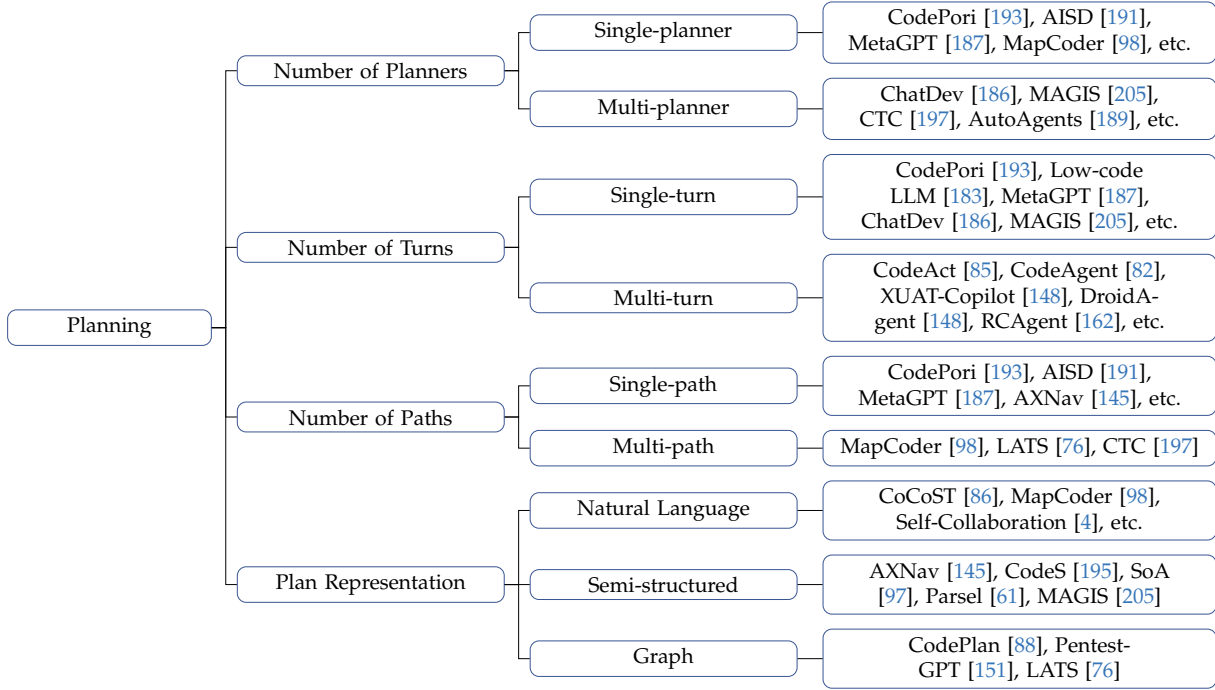


Fig. 15: Taxonomy of Planning Strategies in LLM-based Agents for Software Engineering

5.1.2 Memory

The memory component is a pivotal mechanism responsible for storing the trajectories of historical thoughts, actions, and environmental observations, enabling agents to sustain coherent reasoning and address intricate tasks. In SE, complex development and maintenance tasks generally necessitate agents conducting iterative revisions, wherein historical intermediate information, e.g., generated code and testing reports, significantly impacts integrity and continuity. We then detail the implementation of memory mechanisms in SE from four perspectives: memory duration, ownership, format, and operation. Figure 16 presents the taxonomy of the memory components in existing LLM-based agents for SE.

A. Memory Duration. Inspired by human memory systems, agent memory can be classified into short-term memory and long-term memory based on the memory duration.

A.1: Short-term Memory. Short-term memory, also known as working memory [221], is integrated into agents to enhance their ability to sustain trajectories of the current ongoing task and is frequently used when multi-turn interactions are involved. In SE, there are some predominant patterns of short-term memory.

- *Dialog Records.* This pattern is generally used to memorize the pure dialog history among agents and is typically in the form of history summary [148], [189] and multi-turn instruction-response pairs [186], [197], [198]. It is straightforward to implement and can offer a thorough and detailed historical record of the task-solving process. However, the weakness is that the dialog history can be lengthy and contain irrelevant and redundant information.
- *Action-Observation-Critique Records.* While dialog history concentrates on the thoughts and responses among

agents, some works highlight the interaction between agents and the environment by memorizing the action-observation sequences. Moreover, the critique information is also retained in case certain reflection mechanisms are introduced [144]. This pattern has been adopted in SE tasks that necessitate iterative feedback from the environment, e.g., mobile app testing [139], [144], [148], wherein operations on widgets in each turn should be memorized to facilitate the next-turn decision-making, and iterative code generation [62], [88], [187], [189], wherein the previous editing, execution, or debugging history serves as important information for code revision.

- *Intermediate Outputs.* Some agents store the outputs of previous turns in short-term memory to avoid overrunning the limited space as well as being overly influenced by irrelevant or inaccurate chat history. For example, in E&V [107], only intermediate analysis results are summarized to avoid inconsistency with the previously generated outputs. In SoA [97], to implement a self-organized framework, each agent is equipped with memory that stores the self-generated code and unit tests. These intermediate results allow delayed test execution and code modification for agents in different layers, facilitating hierarchical collaborative code generation.

A.2: Long-term Memory. Long-term memory, on the other hand, is used to memorize valuable experiences of historical tasks, which can be recalled by agents when solving unseen tasks. Due to extensive trajectories, long-term memory commonly uses distilling techniques or only stores the pivotal information.

- *Distilled Trajectory.* The entire task execution trajectory may involve extensive context, and given the limited memory space, it can be challenging to store it all completely. As a result, distilling techniques have been

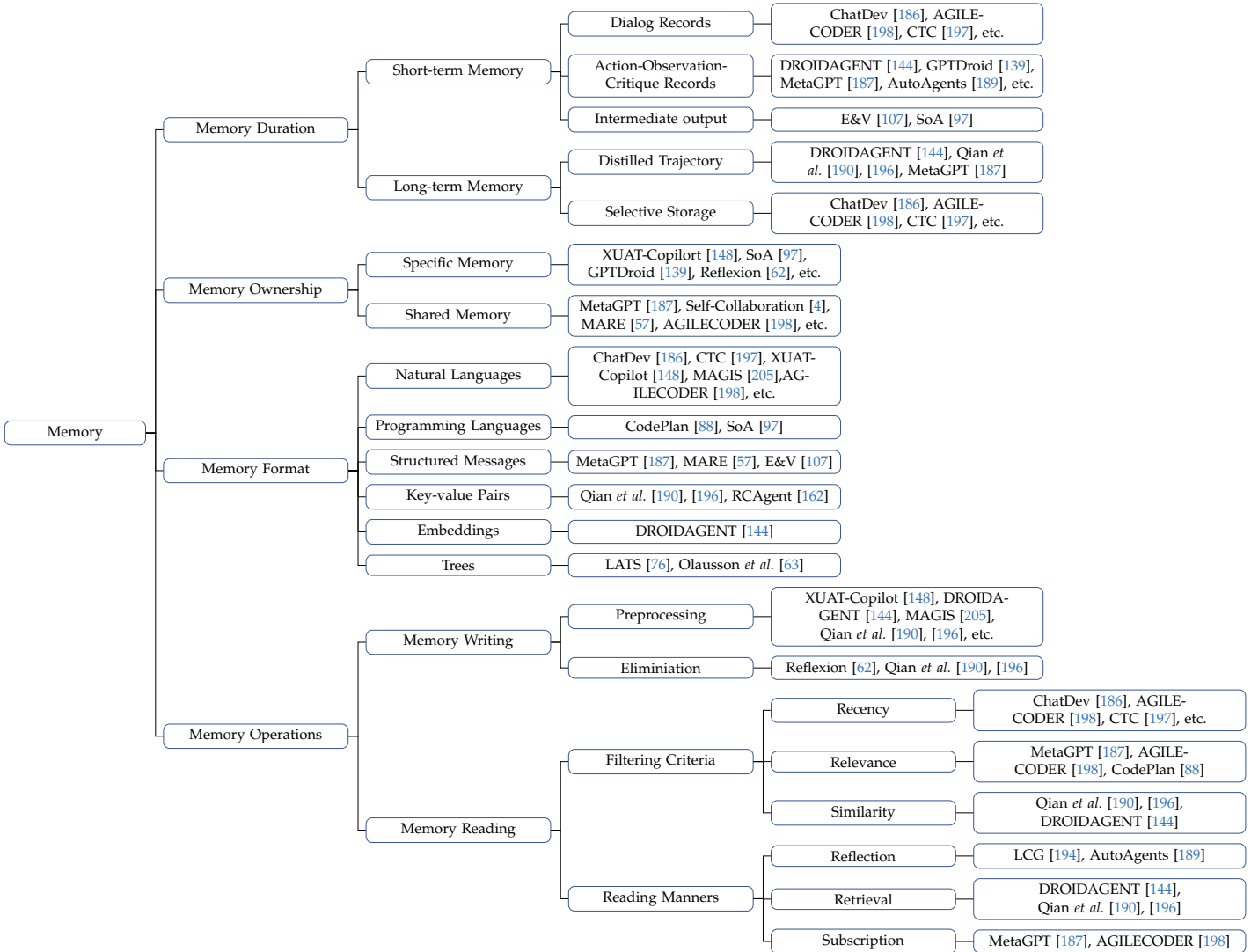


Fig. 16: Taxonomy of Memory Design in LLM-based Agents for SE

proposed, *e.g.*, trajectory summarization [144], [187] and shortcut extraction [190], [196]. These distilled records retain the task execution process in a more concise manner, thereby alleviating the burden on limited memory and prompt windows.

- *Selective Storage*. Another manner to save long-term memory space is to store vital data of each task, *e.g.*, the final results [186], [189], [197], [198], reflections [62], [189], and action-observations [76], [144], [162]. Compared to complete historical trajectories, these data highlight the pivotal trace information, which can still retain the effects and feedback of previous tasks.

B. Memory Ownership. In agent systems, the memory module can be designed to serve specific agents or to serve all agents. Based on its ownership, we categorize memory into Specific Memory and Shared Memory.

B.1: Specific Memory. Specific memory is an agent mechanism designed specifically for a limited group of agents. This type of memory has strict pre-defined usage regulations, only storing and serving specific agents in the workflow [62], [76], [88], [97], [98], [107], [139], [144], [148], [162], [186], [187], [189], [190], [196]–[198], [205]. For example, in

SoA [97], each agent is equipped with individual memory for storing its own generated code fragments and unit tests, which will be used to evaluate the correctness of the final code and provide feedback to the agent for modification.

B.2: Shared Memory. Shared memory, on the other hand, serves all agents by maintaining the record of their outputs and offering essential historical data. In most cases, shared memory serves as a dynamic information exchange hub in the intricate SE environment, which is akin to the traditional blackboard system [222]. Generally, information stored in the shared memory is the intermediate results of previous phases, hence the agents from subsequent phases can obtain necessary information in a more convenient manner [4], [57], [160], [187], [198]. Representative work like MetaGPT [187] introduces a shared message pool, which saves artifacts from different agent roles, *e.g.*, the product requirement documents from the product manager. Another typical application of shared memory is to store comments in a decentralized debate scenario. Specifically, LCG_{scrum} [194] simulates the Sprint Meeting by providing a shared buffer, storing the problem and the discussion comment of all participated agents from which the product manager

could extract a list of user stories.

C. Memory Format. In this section, we elaborate on the format of data stored in the memory. In SE tasks, the most commonly used storage formats include natural languages, program languages, structured messages, key-value pairs, embeddings, and trees.

C.1: Natural Languages. LLM-based agents solve tasks specified in natural language, which is thus the most fundamental and prevalent data format in memory [4], [62], [139], [148], [186], [189], [194], [197], [198], [205]. The advantage of raw natural language is that it allows for a more flexible storage of trajectories, thereby enhancing the universality. Moreover, raw natural language can better preserve the integrity of the original dialogue, which minimizes the loss and distortion of essential information.

C.2: Programming Languages. Some agents directly store the generated code for subsequent utilization [88], [97]. For example, SoA [97] is a hierarchical code generation framework with each agent focusing on single function implementation. It stores the generated function code and unit tests in the memory for testing, modification, and aggregation.

C.3: Structured Messages. In this format, memory is organized as a list of messages with multiple attributes. The strength of this format is that it allows data to be stored in a structured manner, making it more convenient for indexing and processing. Moreover, it can store vital metadata of the message, *e.g.*, message source and destination, task type, etc., so it is easier for agents to trace and subscribe to required messages, making it commonly used in *shared memory*. Representative works like MetaGPT [187] and MARE [57], both wrap the artifacts of each agent as informative messages, involving the original content, instruction, task name, sender, and receiver, etc. In E&V [107], intermediate results of previous turns will be summarized and stored in JSON format.

C.4: Key-value Pairs. In this format, information received from the agents is stored in an external memory, with a key extracted for the agents to query required history memories [162], [190], [196]. More specifically, in Co-Learning [190], shortcuts are extracted from the trajectories to construct two key-value databases: the solution-to-instruction database for the instructor, and the instruction-to-solution database for the assistant. RCAgent [162] stores the whole observation body in a key-value store, remaining a snapshot key for agents for query details.

C.5: Embeddings. In this format, the memory is embedded into a vector, which can help retrieve the most relevant task experiences. Representative work like DROIDAGENT [144] embeds the textual history into vectors and stores them in an external embedding database. Compared to text similarity retrieval, it can further provide semantic similarity retrieval.

C.6: Trees. Some approaches construct a tree or graph for memorizing, especially in scenarios requiring flexible extension or path tracing. For example, in LATS [76], the task-solving process is modeled into a tree with each node representing a state with the instruction, the action, and the observation, and then an extended Monte Carlo Tree Search algorithm can be integrated. Similarly, Olausson *et al.* [63] propose a repair tree that stores multiple generation-feedback-repair paths.

D. Memory Operations. We categorize operations on memory into two main sections: memory writing and memory reading.

D.1: Memory Writing. The purpose of memory writing is to store essential information in the memory. In SE, the most commonly considered problems include how to process the received information (memory preprocessing) and how to avoid memory overflow (memory overflow).

- *Memory Preprocessing.* Information stored in memory is usually the raw task execution trajectories [186], [197], [198]. However, considering that the raw task trajectories might be lengthy, distilling approaches have been proposed to retain a more informative summary in the memory [62], [107], [144], [148], [187], [190], [196], [205]. For instance, in XUAT-Copilot [148], dialog and action history are stored in working memory as summarized texts. Moreover, Co-Learning [190] proposes a novel distilling approach by first constructing a task execution graph and then extracting shortcuts linking non-adjacent solution nodes, which can serve as solution refinement paths for future tasks.
- *Memory Elimination.* The limited memory storage and prompt window size result in finite memory records. When overflow occurs, some records must be forgotten. For example, in Reflexion [62], the past experiences are stored in a sliding window with a maximum number of 3 to avoid exceeding the prompt window. Additionally, low-quality and rarely-used data also consume memory storage space. Previous research [190] sets a threshold to filter out experiences with limited information. Further, an elimination mechanism based on the usage frequency is introduced to exclude rarely-used experiences [196].

D.2: Memory Reading. Memory reading aims at obtaining the required task history and experiences from the memory module. We elaborate on the memory reading process from two perspectives: the filtering criteria required historical records and reading manners to obtain these records.

Filtering Criteria. The factors influencing whether a historical record should be integrated into the current task can be chiefly categorized into three parts: *recency* [88], [139], [144], [148], [186], [187], [197], [198], [205], *relevance* [88], [187], [198] and *similarity* [144], [190], [196]. Agents integrate the most relevant and similar task experiences to provide the best reference for the current task. Additionally, the preference and weight of these factors vary across different works. In DROIDAGENT [144], the planner agent considers the 20 most *recent* task summaries and the 5 most *similar* task knowledge items. In MAGIS [205], the agent uses the *most recent* summary of a code file to identify differences. In MetaGPT [187] and AGILECODER [198], agents retrieve only *relevant* messages from shared memory based on their roles.

Reading Manners. In SE, instead of directly providing the raw memory to the agent [4], [139], [160], [186], [187], [197], [198], [205], researchers prefer using three manners for obtaining relevant memory: *reflection*, *retrieval*, and *subscription*.

- *Reflection.* Reflection refers to extracting pivotal experiences from the extensive trajectory memory [189], [194]. For example, in AutoAgents [189], the dynamic memory mechanism is designed to instruct an agent to extract in-

sights from long-term memory that will serve the current action. In LCG_{scrum} [194], the product manager summarizes the comments collected from all agents and extracts a list of user stories to implement.

- *Retrieval*. In this manner, the memory is retrieved based on its text or semantic similarity with the current tasks [144], [190], [196]. For example, in Co-Learning [190], the reasoning module uses the prompt as a query to retrieve similar shortcuts from the constructed experience pool, which will serve as examples to facilitate future reasoning. In DROIDAGENT [144], the past tasks and widgets with similar GUI state embeddings are retrieved by comparing the cosine similarity.
- *Subscription*. The subscription mechanism is chiefly used in shared memory. It permits agents to directly obtain required information according to their roles, without additional interaction costs with other agents, thus improving efficiency. Representative works include MetaGPT [187] and AGILECODER [198], both adopting this kind of publish-subscribe mechanism.

5.1.3 Perception

Existing LLM-based agents for SE primarily adopt two perception paradigms: textual input perception and visual input perception.

A. Textual Input. Text can flexibly express the intent, information, and knowledge. In SE, the majority of historical data, *e.g.*, documentation, code, and issues, is stored in the textual form. This alignment with the strengths of LLMs in processing natural language makes textual input the predominant form of perception for agents for SE. Textual input in existing agents can be further categorized into natural language input (*i.e.*, instructions and auxiliary information collected from the environment) and programming language input (*i.e.*, the code context). For example, in NL2Code tasks [186], [187], user requirements and function descriptions are provided as the instruction to agents. But in some code-related tasks, *e.g.*, software testing [126], [127] and debugging [159], [160], [171], [174], the target code can also be provided for analysis. Specifically, in repository-level tasks such as issue-resolution [206]–[208], repository-level fault localization [160], and code edits [88], only a portion of code snippets are provided due to context length limitations, with further inspections achieved through navigation in the code repository.

B. Visual Input. Images represent a two-dimensional medium for storing information. In traditional SE scenarios, there is also a portion of data presented in image form, *e.g.*, UML diagrams [223] and UI pages. A small number of existing agents use this visual information to enhance their understanding of target tasks. For example, in MetaGPT [187], the product manager creates a competitive quadrant chart for the architect, who then provides system architecture and sequence flow diagrams to the engineer agents.

Visual input is more widely used in mobile app testing tasks, since the clickable widgets can be located easily in screenshots. Previous work [144] has also discovered that the widgets in some apps might be presented in raw pictures without any textual information [144]. Therefore, in XUAT-Copilot [148] and AXNav [145], the screenshot of the current page is provided to the agent to visualize the accessible

widgets. Just as humans use their eyes to interpret images, these agents integrate external visual models to process and understand image data. XUAT-Copilot [148] uses the SegLink++ model [224] for detecting bounding boxes and a ConvNeXts model [225] for text recognition.

5.1.4 Action

The action component of existing LLM-based agents for SE primarily involves using external tools to extend their capabilities beyond the interactive dialogue typical of standalone LLMs. Figure 17 summarizes the tools used in these agent systems.

A. Searching Tools. In SE, agents frequently use searching tools to retrieve relevant information (*e.g.*, documentation or code snippets) that can be helpful for the task completion.

A.1: Web Searching. Online search engine tools use community and tutorial websites to offer programmers accurate and practical suggestions based on shared experiences and Q&A. When faced with gaps in specific domain knowledge, programmers distill their needs into a query and use existing search engines (*e.g.*, Google, Bing, WikiSearch) to find the necessary information, an approach that can also be employed by agents [78], [96], [105], [153], [187], [193]. For example, some agents [78], [81], [82], [188] use DuckDuckGo [228] to search the relative content such as APIs. Paranjape *et al.* [102] employ SerpAPI [229] and extract answer box snippets when they are available or combine the Top-2 search result snippets together. He *et al.* [230] query Google and then extract pertinent information to construct prompts for LLMs. Depending on the task at hand, the search space can be restricted to specific websites (*e.g.*, as StackOverflow) or certain websites can be blocked to prevent data leakage.

A.2: Knowledge Base Searching. Besides using a web searching tool to externally collect the information from the wide, it is also common for existing agents to retrieve relevant knowledge from the self-established knowledge base (*e.g.*, memory pool or code repository). In particular, similarity and string matching are two common retrieval strategies. Similarity-based retrieval approaches include sparse word-bag and dense text embedding. Both approaches vectorize codes or documents, calculating similarity between the query and the segment in knowledge base to obtain relevant information. The sparse word-bag approach (*e.g.*, BM25 [78], [82]) vectorizes text while partially preserving its semantics. Dense text embedding model such as dual-encoder encodes the text into embedding vectors and calculates their cosine similarity [64], [77], [105], [110], [144], [151], [162], [190]. String matching approaches directly split the given code element name and match it within the knowledge base [107], [171], which is used for locating files within a repository by the file name.

B. File Operation. As SE activities frequently access massive files especially for the code repository and documentation, it is common for agent systems [113], [153], [171], [172], [207], [210] to use file operations including shell commands (*e.g.*, Linux shell) or the code utils (*e.g.*, Python *os* package) for file browsing, file adding, file deleting, and file editing. For example, for file browsing, agents open files

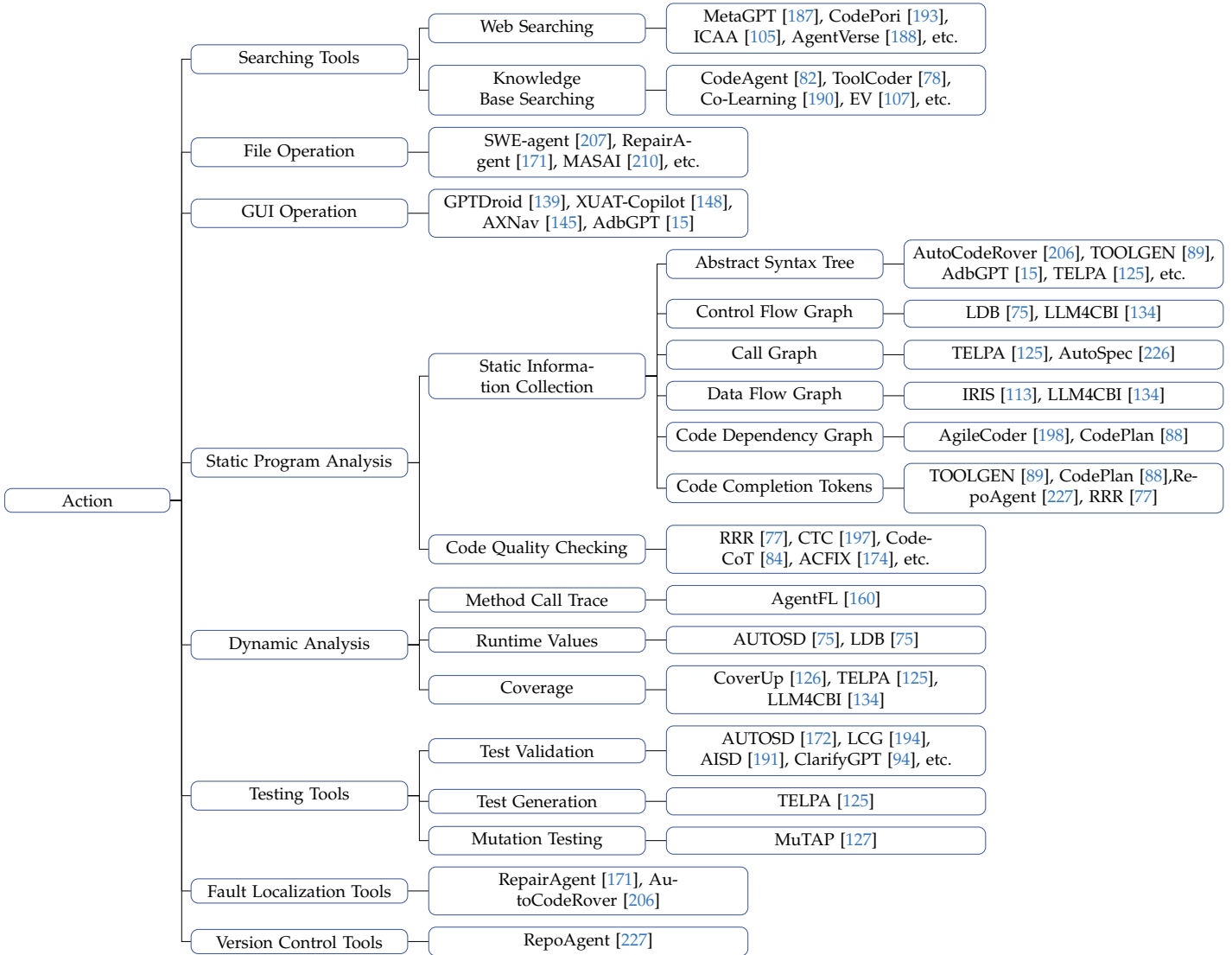


Fig. 17: Taxonomy of Action Components in LLM-based Agents for SE

based on their paths, scroll through the contents, and jump to specific lines.

C. GUI Operation. For SE activities related to software with GUI, it is necessary to enable various GUI interaction operations for agent systems [15], [139], [145], [148], including clicking, text input, scrolling, swiping, returning, and termination. In particular, for UI element identification, they use visual and text recognition models (e.g., SegLink++ [231], Screen Recognition [232], and ConvNeXts [233]), dump (e.g., Android UIAutomator [147]), or parse the UI view hierarchy [15], [139]; then they simulate the testing environment using virtual Android devices (e.g., Genymotion [146], VirtualBox [140], and pyvbox [141]) and autonomously execute or reply actions through tools such as Android Debug Bridge [143] to mimic user interactions. These actions enable agent systems to test in various GUI environments.

D. Static Program Analysis. Static program analysis tools are widely used in agent systems for SE tasks, as they can provide more rigorous code features (e.g., data-flow and control-flow) for LLMs. Existing agents primarily use static

program analysis tools for two purposes: collecting static program information and checking code quality.

D.1: Static Information Collection. Agents invoke static analysis tools to parse the program and to collect additional program information; and the collected information can further help agents better understand the program and thus tackle the relevant tasks. Existing agents primarily collect the following static information.

- *Abstract Syntax Tree (AST).* AST is a common representation to describe the syntactic structure of the source code and is widely used by agents. In particular, the collected ASTs help agents extract syntactic elements (e.g., class names, method names, and variable names) [15], [77], [82], [88], [88], [89], [107], [125], [126], [198], [206], [210], [226] and identify dependency among these code elements [88], [125], [198], [226], [227]. Tree-sitter [161] and ANTLR [234] are AST parsing tools that are widely used in existing agent systems [82], [88], [160], [171], [174], [198], [210].
- *Control Flow Graph (CFG).* LDB [75] uses CFG to divide a program into multiple blocks, making it easier to track intermediate variables with the help of the debug-

ger. LLM4CBI [134] calculates the cyclomatic complexity based on a CFG that represents failed tests and accurately identifies high-complexity blocks of the code. These complicated code blocks would be regarded as the targets for program mutation.

- *Call Graph (CG)*. TELPA [125] constructs a method CG, extracting all call sequences that reach uncovered target methods. Based on these sequences, new test cases are generated to ensure comprehensive coverage of previously untested methods. AutoSpec [226] treats loops as nodes as well, constructing an extended call graph; and it then traverses the CG from the bottom up to generate specifications.
- *Data Flow Graph (DFG)*. IRIS [113] constructs a data flow graph to assist taint analysis, which helps detect security vulnerabilities. LLM4CBI [134] performs data flow analysis to output a list of the most complex variables defined and used in the failed test, which guides test generation.
- *Code Dependency Graph (CDG)*. Agents such as AgileCoder [198] and CodePlan [88] build a Code Dependency Graph for the entire codebase. The graph represents complex relationships between code blocks (e.g., call relationships, inheritance, and import dependencies) and enables the accurate extraction of task-relevant context information (e.g., error trace-back path for repair) within constraints of a limited prompt length. In addition, by dynamically maintaining the CDG, agents can perform incremental analysis in a more efficient way.
- *Code Completion Tokens*. In code generation tasks, it is common for agents [77], [88], [89], [227] to use language servers (e.g., Jedi [235] and EclipseDTLS [236]) to collect candidate tokens at the certain position. In particular, candidate tokens returned by language servers often pass the syntactic violation (e.g., only defined variable names are returned), which can effectively alleviate the hallucinations of standalone LLMs.

D.2: Code Quality Checking. Static analysis tools are also widely used by agent systems to check code quality, e.g., syntactic correctness checking, code format checking, code complexity checking, vulnerability detection, and specifications checking. The checked results can then provide feedback or additional hints for agents to further improve code quality. In particular, existing agents [77], [84], [197] use compilers or interpreter (e.g., GCC or Python) for syntactic correctness checking; existing agents [82], [192] use Black [237] and nuXmv [238] for code format checking; the agent in [134] uses Oclint [135] and srcSlice [136] for code complexity checking; existing agents [134], [174] use static analysis tools such as Frama-C [138] and Slither [239] to detect vulnerabilities; the agent in [226] uses static tools (e.g., Frama-C) to verify the satisfiability and sufficiency of generated specifications.

E. Dynamic Analysis. In addition to static analysis, existing agents also use dynamic analysis tools to collect dynamic runtime information (i.e., method call trace, runtime values, and coverage) that can further provide runtime behaviors for agents.

- *Method Call Trace*. AgentFL [160] uses the `java.lang.instrument` package [240] to record all method call traces during the execution of failed tests, which can

facilitate more accurate fault localization.

- *Runtime Values*. Some agents [75], [172] mimic manual debugging to set breakpoints, so as to capture runtime values for variables. The runtime values can be integrated into the prompt along with requirements to aid in defect localization.
- *Coverage*. Coverage serves as important feedback for whether each code element is executed by tests or not. For example, some agents [125], [126], [134] leverage tools such as SlipCover [129], Pynguin [241], and Gcov [137] to collect the coverage information.

F. Testing Tools. Test cases validate whether the software behaviors violate the specifications, and it is common for agents in SE to invoke testing tools, including tools for test validation, test generation, and mutation testing.

F.1: Test Validation. Validating the software with test execution frameworks (e.g., PyTest, unittest, or JUnit) can reveal the runtime errors and test failures, which are widely used in existing agent systems [62]–[67], [71]–[74], [79], [81]–[85], [87], [91], [92], [94], [95], [97], [98], [102], [122], [123], [125], [128], [133], [153], [159], [171], [172], [175], [187], [188], [191], [194], [198], [206], [210], [230]. The revealed execution violations can further serve as feedback for agents to improve programs; otherwise, the absence of execution violation can serve as a signal for the correctness of programs (e.g., a plausible patch is found for program repair agents).

F.2: Test Generation Tools. Although an LLM itself has promising capabilities of directly generating test code, traditional test generation tools provide complementary benefits as they are good at generating high-coverage tests in a cost-efficient way. For example, some agents [125] use automated test case generation tools (e.g., Pynguin [241]) to generate an initial set of unit test cases.

F.3: Mutation Testing. Some agents [127] use mutation testing tools (e.g., MutPy [242]) to evaluate the sufficiency of test cases, as killing mutants (i.e., exhibiting different behaviors on the mutated program than the original program) indicates the fault detection capabilities of tests. The mutation testing results can further serve as the feedback for agents to iteratively enhance the tests.

G. Fault Localization Tools. Agent systems [171], [206] can invoke traditional fault localization techniques, especially spectrum-based fault localization tools (e.g., GZoltar [243]) to localize suspicious code elements. For example, RepairAgent [171] invokes GZoltar to get the suspiciousness score of each code element (i.e., the probability of being fault).

H. Version Control Tools. Version control systems manage the changes of various files in a repository such as changes in code, configuration files, or documentation throughout the software development process. Some agents [227] involving managing an entire repository often leverage version control tools. For example, RepoAgent [227] uses Git [244] to track changes in code files and promptly synchronize updates to project documentation.

5.2 Multi-agent System

Based on our statistics, 52.8% of existing agents for SE are multi-agent systems. These systems benefit from the division of specialized roles and coordination among agents,

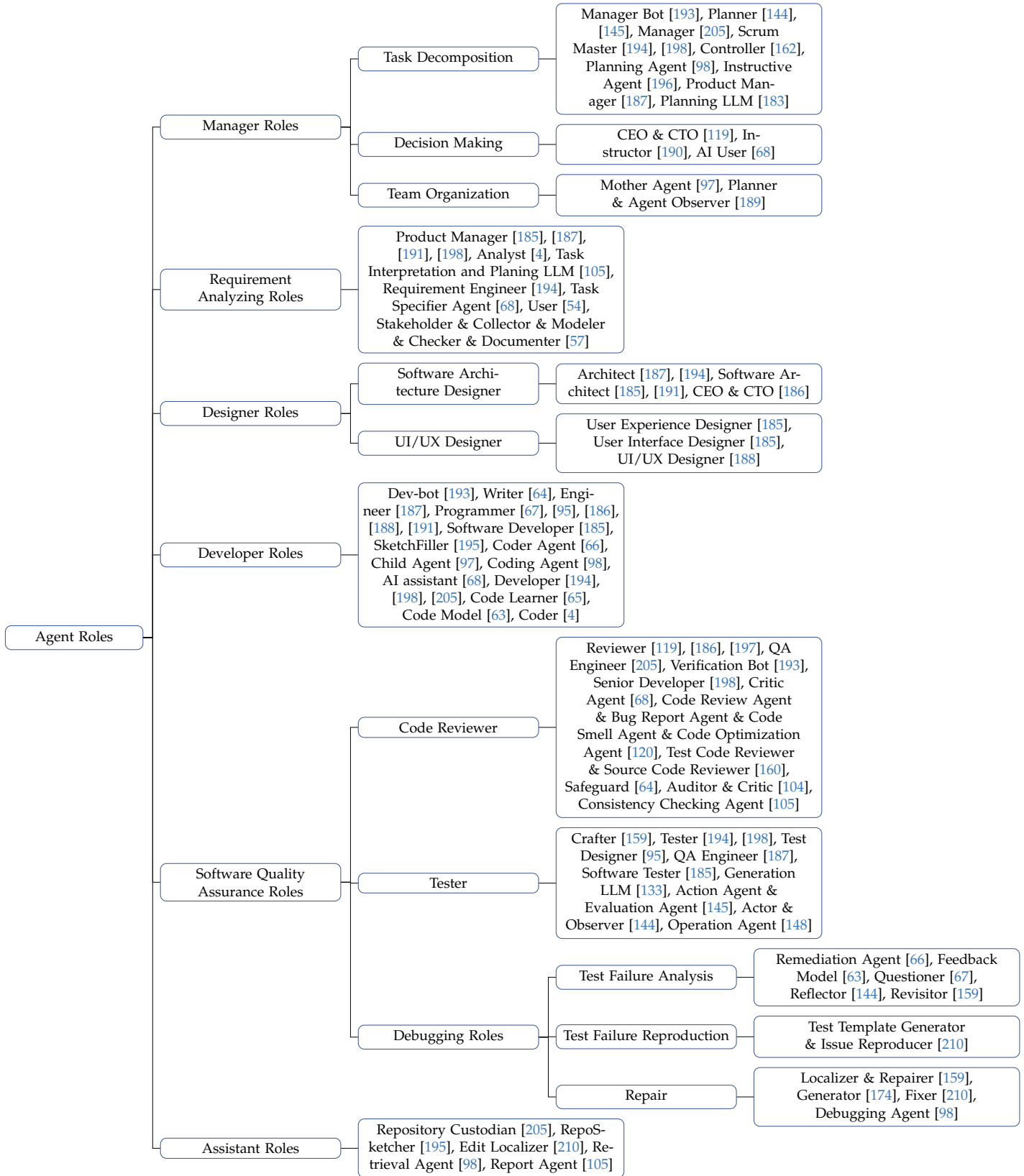


Fig. 18: Taxonomy of Agent Roles in LLM-based Multi-agent Systems for SE

which effectively addresses the complexity of SE tasks, particularly for end-to-end activities spanning multiple phases. This section provides an overview of existing multi-agent systems for SE, with a focus on their agent roles and coordination mechanisms.

5.2.1 Agent Roles

In a multi-agent system, each agent is typically assigned a specialized role designed to address specific tasks. This role assignment forms the foundation of the system, directly influencing task decomposition and agent coordination. Specifically, role assignment mainly delineates duties, available actions, attributes, and constraints of roles. It makes agents as experts of corresponding tasks. Figure 18 summarizes common agent roles in existing multi-agent systems for SE.

A: Manager Roles. Managerial roles, such as CEO, CTO, commander, and controller, serve as the leaders of a multi-agent team. These roles are responsible for making decisions, planning, task decomposition and assignment, and overseeing team coordination.

A.1: Roles for Task Decomposition. To enhance the overall system performance, managers break down a project into manageable sub-tasks and draw up a guiding plan for developers or testers to execute [98], [144], [145], [162], [193], [196], [198], [205]. They analyze problem statements, facilitate discussions on issues among various agent roles [194], review design documents submitted by designers [187], and incorporate related information gathered by assistants. Subsequently, they produce a specific task list or implementation blueprint, which may be presented in either natural language or as a structured workflow [183].

A.2: Roles for Decision Making. This role is designed for orchestrating collaboration within a team and providing further guidance for task execution. For example, CEO and CTO in CodeAgent [119] communicate with staff and make high-level decisions. Similarly, the instructor in the agent [190] and the AI user within the CAMEL system [68] issue directions to working agents.

A.3: Roles for Team Organization. This role is primarily designed for flexibly deciding the constitution of the agent team, *i.e.*, what roles are included in the team. The main benefits of including such roles are to flexibly optimize costs and better meet project demands. For instance, SoA [97] sets the mother agent, which generates new mother or child agents and designates concrete tasks (*e.g.*, unimplemented functions) to them. AutoAgents [189] includes a planner agent and an observer agent, which collaborate to assemble a team for particular tasks. The planner agent is responsible for assigning existing LLM agent roles or generating new ones, while the observer agent assesses and reviews the relevant roles. These roles are represented in a structured JSON format, encapsulating details such as name, description, available tools, suggestions, and prompts to guide agent behaviors.

B. Requirement Analyzing Roles. These roles are primarily responsible for analyzing software requirements, such as translating vague and preliminary user concepts into a coherent and structured format. Existing agents [4], [105], [185], [187], [194], [198] include such roles (*e.g.*, product manager [191] or task specifier [68]) to identify key

requirement elements and intended objectives for a precise and organized requirement document, which may range from an elaborated task or function description [68] to a formal software requirement specification [187].

In addition, some agents further design more fine-grained roles for requirements analysis. For example, Elicitor [54] incorporates a set of *User* agents to identify diverse user requirements by mimicking user perspectives and conducting interviews for exploring potential user needs; MARE [57] uses a requirements engineering team (*i.e.*, stakeholder, collector, modeler, checker, and documenter) to produce requirements specifications. The requirements engineering process is segmented into four sub-tasks corresponding to specific roles, seamlessly transitioning rough user ideas to precise requirement specifications.

C. Designer Roles. Designer roles take input information on requirements (such as detailed task descriptions and use cases) and shape the software architecture and system integration.

C.1: Software Architecture Designer. This role is responsible for conceptualizing and defining the high-level structure of software, *e.g.*, the *software architect* role in agents [185]–[187], [191], [194]. They create a design document that serves as a blueprint for the subsequent stages of development; and the design document can be presented in various forms, including natural language descriptions, structured formats (*e.g.*, JSON for listing project architecture files), and graphical representations (*e.g.*, class diagrams and sequence flowcharts [187]).

C.2: UI/UX Designer. This role primarily focuses on crafting the visual and interactive aspects of the software interface, such as user experience (UX) designer or user interface (UI) designer roles in agents [185], [188].

D: Developer Roles. Developers take a vital role in software development and maintenance activities, which is one of the most common roles (*e.g.*, also called as programmer or coder) in existing agents [64], [66], [68], [97], [98], [185]–[188], [193], [195] for tasks involving code generation. In accordance with software design schemes, task plans provided by other agents, or user requirements, the developer roles generate or finalize code at various levels (*i.e.*, from function to file and even project levels). In addition, the developer roles also engage in the code refinement process, which refines their previously generated code [4], [63], [65], [67], [95], [194], [205]. Furthermore, the developer roles can be set to meet more customized standards, such as elucidating their work through supplementary docstrings or adhering to particular coding criteria [191], [198].

E: Software Quality Assurance Roles. Agent systems include roles dedicated to software quality assurance, similar to real-world QA teams. These roles typically encompass code reviewers, testers, and debuggers, each focused on checking and improving software quality.

E.1: Code Reviewer. This role is responsible for identifying potential software quality issues by statically inspecting the software without execution. For example, some agents [119], [186], [188], [193], [197], [205] include such roles to review generated code or patches; AgileCoder [198] and CAMEL [68] include the roles such as senior developer or critic agent to offer suggestions for enhancement; the agent in [120] sets up code review agent, bug report agent, code

smell agent, and code optimization agent to access code quality from different aspects; AGENTFL [160] sets test code reviewer and source code reviewer to summarize code behaviour to help fault location; in addition, some agents [64], [104], [105] include such roles (e.g., the auditor agent and the critic agent in GPTLENS [104] and the consistency checking agent in [105]) to detect the vulnerability or implementation issues.

E.2: Tester. The tester roles are widely incorporated in multi-agent systems [95], [133], [144], [145], [148], [159], [185], [187], [194], [198] for software quality assurance, which are mainly responsible for writing new tests or generating testing action sequences. For example, the tester agent in multi-agent systems [95], [159], [194] generates test cases based on relevant code skeleton or patches, requirement documents, existing tests, or rationale for the test; the tester agent in multi-agent systems [144], [148] generates operational actions for GUI tests or systematic functional tests based on the specified requirements.

E.3: Debugging Roles. In multi-agent systems, debugging roles are responsible for diagnosing test failures or unexpected software behaviors.

- *Test Failure Analysis.* Some multi-agent systems include debugging roles in analyzing test reports. For example, the remediation agent in TGen [66] and the feedback module in Self-repair [63] are similarly designed to analyze the test failure reports and relevant faulty code to provide explanations and suggestions; the questioner agent in AutoCoder [67] describes execution errors to help modify the generated code; the reflector agent in DROIDAGENT [144] reflects and summaries on the test results; and some agents [65], [159] further include debugging roles (e.g., the revisitor) to provide explanations for test failures.
- *Test Failure Reproduction.* MASAI [210] uses the test template generator to produce test templates based on the repository information, which further helps the issue reproducer reproduce behaviors described in the issue reports.
- *Repair.* Some multi-agent systems include such roles in fault localization and program repair. For example, agents [98], [159], [174], [210] include roles such as repairer and fixer to generate patches for bugs;

F: Assistant Roles. Assistant roles primarily provide assistance for other agents. For example, the repository custodian in MAGIS [205], the RepoSketcher in CodeS [195], and the edit localizer in MASAI [210] are designed to enhance the comprehension of the target repository architecture for the team; in addition, MapCoder [98] uses the retrieval agent to facilitate memory recall; ICAA [105] introduces the report agent to convert natural language responses into formatted bug reports.

5.2.2 Collaboration Mechanism

The collaboration mechanism is essential for multi-agent systems, which can significantly impact the effectiveness and costs of the entire system. In particular, the collaborative mechanisms of existing multi-agent systems for SE tasks can be categorized into four types: layered structure, circular structure, tree-like structure, and star-like structure. Figure 19 illustrates each structure.

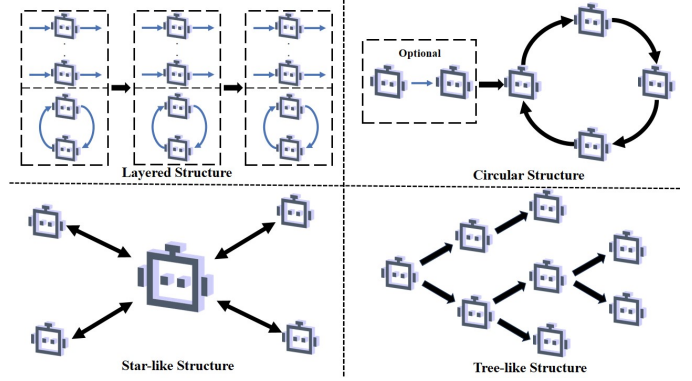


Fig. 19: Multi-agent System Collaboration Mechanisms

A. Layered Structure. It is a hierarchical structure, where tasks are decomposed into several sub-stages and each is assigned to a specific agent or a group of agents selected from the agent pool. Agents between different stages collaborate in a sequential manner, i.e., they receive intermediate results from agents in the previous stage as input and produce their processed data to agents in the next stage. For example, the workflow within agents of [105], [133], [160], [183], [191], [206] is a simple chain, where each agent focuses on its own sub-task and only interacts with adjacent agents. In addition, agents can also refer to the message produced by the previous non-adjacent agents [98], [187], [195]. In the sequential workflow, each sub-task can also be handled by a group of agents [57], [91]. In [119], [186], [197], each sub-task is solved by the conversation between two agent roles. LCG [194] and AgileCoder [198] incorporate even more agents in a single stage. In addition to interactive collaboration, another scenario involves agents within the same layer working in parallel to offer their solutions. These solutions are then combined and passed down to the next layer. For example, GPTLENS [104] employs several auditors to present possible vulnerable functions individually in the generation stage. The work [69] incorporates the majority voting mechanism. DyLAN [70] formulates the LLM-agent collaboration structure into a multi-layered feed-forward network.

B. Circular Structure. This structure typically manifests as multi-turn dialogues or integrates the feedback mechanism within the overall collaborative processes among agents. The feedback loop facilitates the refinement of the agents' outputs through iterative cycles. The circular structure may consist of dual roles in a dialogue or extend to multiple roles, which results in a more complex iterative circle.

B.1: Dual Roles. On the one hand, some agents [63], [66], [67], [95], [174] implement a generation-validation style loop between two agents. In this setup, one agent is tasked with the primary function, such as generating code snippets or patches, while the other agent provides validation feedback, including static analysis results, test outcomes, and improvement suggestions. In the INTERVENOR framework [65], the code learner initiates the process by generating the initial code and subsequently engages in iterative repairs guided by the suggestions from the code teacher. On the other hand, some agents [68], [111], [190] adopt a dual-

agent dialogue cooperative model to achieve objectives. The work [111] facilitates agreement through a discussion between the tester and the developer; and the agents [68], [190] progress through tasks in a step-wise manner, utilizing an instructor-assistant conversational approach.

B.2: Multiple Roles. When more agents are included, the collaborative loop will become more flexible. For example, some works [62], [145], [193] incorporate multiple agents in a larger feedback loop, further dismantling the tasks. DroidAgent [144] embeds an inner loop between the *Actor* and *Observer* in the overall loop between the *Planner* and the *Reflector*.

C. Tree-like Structure. Different from the layered structure, agents in the same layer of the tree-like structure do not cooperate with each other for the same sub-task but focus on their own work. For example, in SoA [97], the mother agent can dynamically spawn new mother or child agents for code generation, thereby forming a tree-like collaboration structure. In MASAI [210], the tests for reproducing issues and the possible patches are generated parallelly and finally aggregated to the ranker to select the best one.

D. Star-like Structure. This structure is a centralized structure, where a central agent serves as the pivot to interact with other agents. For example, the controller agent in the RCAGENT [162] framework can invoke other expert agents as a kind of tool when necessary. The commander in AutoGen [64] coordinates with the writer and the safeguard separately, to craft code and ensure safety. XUAT-Copilot [148] adopts the operation agent as the core, to receive the judgment from the inspection agent and invoke the parameter selection agent to help the action planning.

5.3 Human-Agent Collaboration

While most agents aim to achieve maximum automation, where users only need to propose a request and wait for the agents to complete the task, previous studies [183], [191] show that LLM-based agents often encounter bottlenecks during the software development process. Therefore, some agents incorporate the human-agent cooperation paradigm to further align and enhance the agent performance with human preference and expertise. As summarized in Figure 20, existing agents primarily include the human participation in four phases: planning, requirements, development, and evaluation.

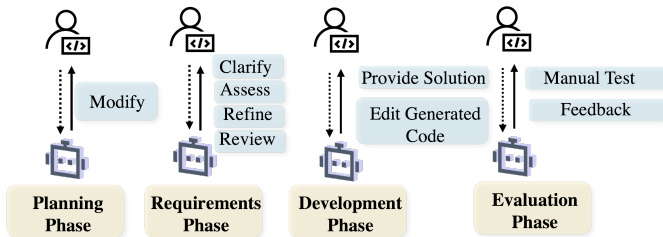


Fig. 20: Human-Agent Collaboration in SE

A. Planning Phase. Some agents include human intervention into the planning stage of the agent workflow. For example, the low-code LLM platform [183] offers users a selection of predefined actions to modify auto-generated workflows. The generated workflows can be checked and

revised by users before execution. However, revising the system design requires a certain level of expertise, so it is optional in some agents such as AISD [191] and LLM4PLC [192].

B. Requirements Phase. The initial requirements (*i.e.*, the task description) provided by users can be ambiguous, which can lead to a gap between the final outputs of the agent system and the user intention. Therefore, it is common for agent systems to further include manual refinement for the requirements. For example, ClarifyGPT [94] and CodeAct [85] employ a dialogue-based interaction with humans; similarly, Sapper IDE [184] and MARE [57] use human feedback to refine the requirements. AISD [191] enables users to assess and refine the generated use cases. In addition, HARA [245] produces a concise and readable summary table of the generated requirements for further manual expert review.

C. Development Phase. Human involvement can be included in the software development phase to guide agents to strategize solutions and overcome potential failures. Flows [91] leverages human-crafted solutions to produce better code for competitive programming challenges. Both AutoGen [64] and LLM4PLC [192] allow users to furnish feedback when necessary, directing the workflow as needed. In CodeS [195], the repository is constructed from a tripartite sketch, assigning users the flexibility to edit each individual layer.

D. Evaluation Phase. Human participation also serves as a post-evaluation mechanism for the outcomes produced by the agent system, which can further ensure the outputs are aligned with user intention. For example, AISD [191] and Prompt Sapper [184] both include human intervention at the acceptance testing phase. Users can conduct manual testing of the final system and the test reports help with the necessary refinements. Similarly, in ART [102], users can enhance agent performance on particular tasks by offering feedback through the modification of the task and tool libraries.

6 RESEARCH OPPORTUNITIES

This section discusses promising research directions and open problems in LLM-based agents for SE.

Evaluation of Agents for SE. Given the emergence of LLM-based agents for SE, it is crucial to develop comprehensive and rigorous evaluation frameworks, including (i) designing more diverse metrics and (ii) constructing higher-quality, more realistic benchmarks.

Metrics. Current evaluations of SE agents primarily focus on their ability to solve specific tasks, such as measuring the success rate of agents on benchmarks such as SWE-bench. However, these evaluations often concentrate on the final success rate without delving into the intermediate states during the agent’s workflow. This lack of fine-grained metrics makes it difficult to assess why agents fail in certain tasks or to what extent they fall short. Given the complexity of SE tasks, failures are common, and without deeper analysis, improving agent performance becomes challenging. Therefore, the design of fine-grained metrics is necessary, allowing researchers to move beyond “black-box” evaluations

and gain insights into the agent’s decision-making process and failure points.

Additionally, existing metrics heavily emphasize effectiveness, leaving trustworthy requirements such as robustness, security, and fairness underexplored. Given the flexibility and autonomy of LLM-based agents, they may exhibit unstable behavior, which can limit their practical application in real-world SE environments. Evaluating these attributes is essential for building trust in these systems.

Another critical consideration is the cost associated with these agents, particularly as they often involve lengthy workflows, frequent LLM invocations, and the management of large datasets. According to our analysis, only 44.3% of the papers we surveyed have explicitly considered the efficiency of agents in SE tasks, incorporating quantitative analyses of time, token consumption, monetary cost, and feedback loops (e.g., tool invocation frequency or inter-agent discussion frequency). These efficiency and computational costs are particularly important when applying agents to large-scale code repositories, complex documentation, or intricate workflows.

Benchmarks. LLM-based agents significantly extend the capabilities of standalone LLMs, showing great promise in tackling more complex, end-to-end SE tasks. However, existing benchmarks used for evaluating these agents often suffer from quality issues. For instance, prior research [211], [215] has identified that the SWE-bench benchmark includes tasks with vague or incomplete issue descriptions, reducing their relevance and applicability.

Moreover, many tasks in current benchmarks are far simpler than real-world SE challenges. As outlined in Table 11, the software generated by LLM-based agents for end-to-end development tends to be relatively small in scale (e.g., consisting of a single function or a few files), which is not representative of the complexity of real-world software projects. In addition, previous work [215] reveals that the majority of tasks (77.8%) in the SWE-bench benchmark can be completed within an hour by an experienced software engineer, further highlighting the gap between benchmark tasks and real-world SE challenges.

To address these shortcomings, future research can focus on creating more realistic, high-quality benchmarks that better reflect the complexity and demands of real-world SE. These improved benchmarks will enable more accurate and meaningful evaluations of LLM-based agents’ capabilities and potential.

Human-Agent Collaboration. Software development is inherently a creative process, transforming human requirements into executable software. As such, aligning agent systems with human preferences and intentions is a critical goal. While some existing agents incorporate human participation at various stages of the workflow (as discussed in Section 5.3), there has been limited exploration of how to more thoroughly integrate human involvement throughout the entire software development life cycle. Additionally, the interaction mechanisms between agents and humans remain underexplored.

Currently, agents mainly involve humans in tasks such as requirements clarification, planning adjustments, coding assistance, or evaluation. However, extending human participation to other phases, such as architecture design,

test generation, code review, and the end-to-end software maintenance process, remains largely unexplored. A deeper integration of human input across these phases could significantly enhance both the quality and adaptability of the agent’s output.

Moreover, designing effective interaction mechanisms is essential for human-agent collaboration. This includes creating user-friendly interfaces for (i) displaying relevant information, such as intermediate outputs from agents, and (ii) collecting user feedback in a streamlined way. Given the complexity of information produced during the agent’s workflow, designing such interfaces presents challenges. For instance, when agents are tasked with generating or maintaining a software repository, simply presenting all code files in a flat format would be resource-intensive and inefficient. Therefore, more sophisticated methods of organizing and representing complex data are required to facilitate effective human-agent interaction.

Perception Modality. Most agents applied to SE tasks primarily rely on textual or visual perception. This is largely because software development and maintenance activities are heavily associated with processing large volumes of code, documentation, and images. However, there is still significant potential to explore and incorporate more diverse perception modalities into these agents.

For example, in the context of programming assistance, most LLM-powered coding agents predominantly use textual input, such as chat interfaces or integrated development environment (IDE) code contexts. Alternative input formats, such as voice commands or user gestures, remain underutilized. Expanding the range of perception modalities could significantly enhance the flexibility and accessibility of coding assistants, allowing users to interact with agents in ways that better suit their individual workflows and preferences.

Furthermore, exploring diverse perception modalities may shape the future of software development and maintenance, offering new opportunities to streamline interactions and improve the efficiency of agent-driven processes.

Applying Agents for More SE Tasks. While existing agents have been deployed across various software SE tasks, several critical phases remain underexplored. As highlighted by our analysis in Section 4, there is a lack of LLM-based agents specifically designed for tasks such as design, verification, and feature maintenance during software development and maintenance.

Developing agent systems tailored to these phases presents unique challenges. Tasks like design and verification require advanced reasoning and comprehension capabilities from the LLM-based agents, extending beyond basic code generation. These tasks demand a deeper understanding of architecture, system logic, and the ability to make informed decisions—skills that traditional LLM-controlled agents may not yet fully possess.

Training Software-oriented LLMs for SE Agents. LLMs are the central component controlling the “brain” of agent systems. Most existing agents for SE rely on LLMs trained on general-purpose data (e.g., ChatGPT [246]) or code-specific data (e.g., Deepseek-Coder [247] and StarCoder [248]). While massive code from GitHub has been leveraged to train LLMs for code, addressing complex SE tasks requires more specialized data. The reason is that

software is not just about code. For example, valuable data from the whole software development life cycle, such as design, architecture, developer discussions/communications, historical code changes, and even dynamic runtime information, remain largely untapped. Incorporating such data into training could lead to the development of more powerful LLMs for software (not just for code), better suited for the unique demands of SE. These enhanced models could form the foundation for more advanced and capable agent systems designed to tackle a wider range of SE tasks.

SE Expertise in Building Agents. Incorporating well-established SE expertise into the design of agent systems is crucial. For instance, widely adopted SE techniques can be integrated as tools or sub-components of agent systems. As discussed in Section 5.1.4, some existing agents already leverage SE toolkits and techniques, but many other SE tools and techniques—such as advanced debugging and testing methods—remain underutilized. Further efforts are needed to comprehensively integrate these tools and techniques into agent systems to enhance their functionality.

In addition, SE domain knowledge can guide the workflow of agent systems. As noted in Section 4.6, some agents for end-to-end software development follow traditional software process models, such as the waterfall or agile models. However, many other software process models remain unexplored. Rather than granting agents full autonomy, existing software development and maintenance methodologies can be used to partially control their workflows. For example, as revealed by the recent Agentless study [211] and also further confirmed by OpenAI [215], LLMs using a simplistic workflow based on traditional fault localization and program repair pipelines can even outperform other more complex, fully autonomous agents. This suggests that leveraging domain expertise from SE can potentially help improve the effectiveness, robustness, efficiency, interpretability, and replicability of agentic solutions.

7 CONCLUSION

In this paper, we have presented a comprehensive and systematic survey of 106 papers on LLM-based agents for SE. We analyzed the current research from both the SE and agent perspectives. From the SE perspective, we analyzed how LLM-based agents are applied across different software development and maintenance activities. From the agent perspective, we focus on the design of components in LLM-based agents for SE. In addition, we discussed open challenges and future directions in this critical domain.

REFERENCES

- [1] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models. *CoRR*, abs/2303.18223, 2023.
- [2] Xinyi Hou, Yanjie Zhao, Yue Liu, Zhou Yang, Kailong Wang, Li Li, Xiapu Luo, David Lo, John C. Grundy, and Haoyu Wang. Large language models for software engineering: A systematic literature review. *CoRR*, abs/2308.10620, 2023.
- [3] Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and Jie M. Zhang. 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*, pages 31–53. IEEE, 2023.
- [4] Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. Self-collaboration code generation via chatgpt. *CoRR*, abs/2304.07590, 2023.
- [5] Burak Yetistiren, Isik Özsoy, Miray Ayerdem, and Eray Tüzün. Evaluating the code quality of ai-assisted code generation tools: An empirical study on github copilot, amazon codewhisperer, and chatgpt. *CoRR*, abs/2304.10778, 2023.
- [6] Jia Li, Yunfei Zhao, Yongmin Li, Ge Li, and Zhi Jin. Towards enhancing in-context learning for code generation. *CoRR*, abs/2303.17780, 2023.
- [7] Junwei Liu, Yixuan Chen, Mingwei Liu, Xin Peng, and Yiling Lou. STALL+: boosting llm-based repository-level code completion with static analysis. *CoRR*, abs/2406.10018, 2024.
- [8] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems*, 36, 2024.
- [9] Yinlin Deng, Chunqiu Steven Xia, Haoran Peng, Chenyuan Yang, and Lingming Zhang. Large language models are zero-shot fuzzers: Fuzzing deep-learning libraries via large language models, 2023.
- [10] Junjie Wang, Yuchao Huang, Chunyang Chen, Zhe Liu, Song Wang, and Qing Wang. Software testing with large language models: Survey, landscape, and vision. *IEEE Trans. Software Eng.*, 50(4):911–936, 2024.
- [11] Caroline Lemieux, Jeevana Priya Inala, Shuvendu K. Lahiri, and Siddhartha Sen. Codamosa: Escaping coverage plateaus in test generation with pre-trained large language models. In *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*, pages 919–931. IEEE, 2023.
- [12] Chunqiu Steven Xia and Lingming Zhang. Less training, more repairing please: revisiting automated program repair via zero-shot learning. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 959–971, 2022.
- [13] Sungmin Kang, Gabin An, and Shin Yoo. A preliminary evaluation of llm-based fault localization. *CoRR*, abs/2308.05487, 2023.
- [14] Harshit Joshi, José Pablo Cambronero Sánchez, Sumit Gulwani, Vu Le, Gust Verbruggen, and Ivan Radicek. Repair is nearly generation: Multilingual program repair with llms. In Brian Williams, Yiling Chen, and Jennifer Neville, editors, *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 5131–5140. AAAI Press, 2023.
- [15] Sidong Feng and Chunyang Chen. Prompting is all you need: Automated android bug replay with large language models. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering, ICSE 2024, Lisbon, Portugal, April 14-20, 2024*, pages 67:1–67:13. ACM, 2024.
- [16] Chunqiu Steven Xia, Yuxiang Wei, and Lingming Zhang. Automated program repair in the era of large pre-trained language models. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pages 1482–1494. IEEE, 2023.
- [17] Nan Jiang, Kevin Liu, Thibaud Lutellier, and Lin Tan. Impact of code language models on automated program repair. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pages 1430–1442. IEEE, 2023.
- [18] Alexander Shypula, Aman Madaan, Yimeng Zeng, Uri Alon, Jacob R. Gardner, Yiming Yang, Milad Hashemi, Graham Neubig, Parthasarathy Ranganathan, Osbert Bastani, and Amir Yazdanbakhsh. Learning performance-improving code edits. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- [19] Russell A Poldrack, Thomas Lu, and Gasper Begus. Ai-assisted coding: Experiments with GPT-4. *CoRR*, abs/2304.13187, 2023.
- [20] Chris Cummins, Volker Seeker, Dejan Grubisic, Mostafa El-houshi, Youwei Liang, Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Kim M. Hazelwood, Gabriel Synnaeve, and Hugh Leather. Large language models for compiler optimization. *CoRR*, abs/2309.07062, 2023.

- [21] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. The rise and potential of large language model based agents: A survey. *CoRR*, abs/2309.07864, 2023.
- [22] CHCR Ribeiro. Reinforcement learning agents. *Artificial intelligence review*, 17:223–250, 2002.
- [23] Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A survey. *Journal of artificial intelligence research*, 4:237–285, 1996.
- [24] Marvin Minsky. Steps toward artificial intelligence. *Proceedings of the IRE*, 49(1):8–30, 1961.
- [25] Charles Isbell, Christian R Shelton, Michael Kearns, Satinder Singh, and Peter Stone. A social reinforcement learning agent. In *Proceedings of the fifth international conference on Autonomous agents*, pages 377–384, 2001.
- [26] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Jirong Wen. A survey on large language model based autonomous agents. *Frontiers Comput. Sci.*, 18(6):186345, 2024.
- [27] Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. A survey on the memory mechanism of large language model based agents. *CoRR*, abs/2404.13501, 2024.
- [28] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. *CoRR*, abs/2402.01680, 2024.
- [29] Yuheng Cheng, Ceyao Zhang, Zhengwen Zhang, Xiangrui Meng, Sirui Hong, Wenhao Li, Zihao Wang, Zekai Wang, Feng Yin, Junhua Zhao, and Xiuqiang He. Exploring large language model based intelligent agents: Definitions, methods, and prospects. *CoRR*, abs/2401.03428, 2024.
- [30] Grégoire Mialon, Roberto Dessi, Maria Lomeli, Christoforos Nalmpantis, Ramakanth Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave, Yann LeCun, and Thomas Scialom. Augmented language models: a survey. *Trans. Mach. Learn. Res.*, 2023, 2023.
- [31] Yue Liu, Sin Kit Lo, Qinghua Lu, Liming Zhu, Dehai Zhao, Xiwei Xu, Stefan Harrer, and Jon Whittle. Agent design pattern catalogue: A collection of architectural patterns for foundation model based agents. *CoRR*, abs/2405.10467, 2024.
- [32] Saikat Barua. Exploring autonomous agents through the lens of large language models: A review. *CoRR*, abs/2404.04442, 2024.
- [33] Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. A survey on large language models for code generation. *CoRR*, abs/2406.00515, 2024.
- [34] Junda He, Christoph Treude, and David Lo. Llm-based multi-agent systems for software engineering: Vision and the road ahead. *arXiv preprint arXiv:2404.04834*, 2024.
- [35] Zhenpeng Chen, Jie M. Zhang, Max Hort, Mark Harman, and Federica Sarro. Fairness testing: A comprehensive survey and analysis of trends. *ACM Trans. Softw. Eng. Methodol.*, 33(5):137:1–137:59, 2024.
- [36] Junjie Chen, Jibesh Patra, Michael Pradel, Yingfei Xiong, Hongyu Zhang, Dan Hao, and Lu Zhang. A survey of compiler testing. *ACM Computing Surveys*, 53(1):4:1–4:36, 2020.
- [37] George Mathew, Amritanshu Agrawal, and Tim Menzies. Finding trends in software research. *IEEE Trans. Software Eng.*, 49(4):1397–1410, 2023.
- [38] Li Zhang, Jia-Hao Tian, Jing Jiang, Yi-Jun Liu, Meng-Yuan Pu, and Tao Yue. Empirical research in software engineering - A literature survey. *J. Comput. Sci. Technol.*, 33(5):876–899, 2018.
- [39] Jie M. Zhang, Mark Harman, Lei Ma, and Yang Liu. Machine learning testing: Survey, landscapes and horizons. *IEEE Transactions on Software Engineering*, 48(2):1–36, 2022.
- [40] DBLP. <https://dblp.org>, 2024.
- [41] 7 million publications. <https://blog.dblp.org/2024/01/01/7-million-publications/>, 2024.
- [42] arXiv. <https://arxiv.org/>, 2024.
- [43] Bin Lin, Nathan Cassee, Alexander Serebrenik, Gabriele Bavota, Nicole Novielli, and Michele Lanza. Opinion mining for software development: A systematic literature review. *ACM Trans. Softw. Eng. Methodol.*, 31(3):38:1–38:41, 2022.
- [44] Klaus Pohl. *Requirements engineering: An overview*. Citeseer, 1996.
- [45] Bashar Nuseibeh and Steve Easterbrook. Requirements engineering: a roadmap. In *Proceedings of the Conference on the Future of Software Engineering*, pages 35–46, 2000.
- [46] Vivek Shukla, Dharendra Pandey, and Raj Shree. Requirements engineering: a survey. *Communications on Applied Electronics*, 3(5):28–31, 2015.
- [47] Grady Booch, Ivar Jacobson, James Rumbaugh, et al. The unified modeling language. *Unix Review*, 14(13):5, 1996.
- [48] Michael Johnson, Robert Rosebrugh, and RJ Wood. Entity-relationship-attribute designs and sketches. *Theory and Applications of Categories*, 10(3):94–112, 2002.
- [49] Xianchang Luo, Yinxing Xue, Zhenchang Xing, and Jiamou Sun. PRCBERT: prompt learning for requirement classification using bert-based pretrained language models. In *37th IEEE/ACM International Conference on Automated Software Engineering, ASE 2022, Rochester, MI, USA, October 10-14, 2022*, pages 75:1–75:13. ACM, 2022.
- [50] Madhava Krishna, Bhagesh Gaur, Arsh Verma, and Pankaj Jalote. Using llms in software requirements specifications: An empirical evaluation. *arXiv preprint arXiv:2404.17842*, 2024.
- [51] Jianzhang Zhang, Yiyang Chen, Chuang Liu, Nan Niu, and Yinglin Wang. Empirical evaluation of chatgpt on requirements information retrieval under zero-shot setting. In *2023 International Conference on Intelligent Computing and Next Generation Networks (ICNGN)*, pages 1–6. IEEE, 2023.
- [52] Krishna Ronanki, Beatriz Cabrero-Daniel, and Christian Berger. Chatgpt as a tool for user story quality evaluation: Trustworthy out of the box? In *International Conference on Agile Software Development*, pages 173–181. Springer, 2022.
- [53] Dipeeka Luitel, Shabnam Hassani, and Mehrdad Sabetzadeh. Improving requirements completeness: Automated assistance through large language models. *Requirements Engineering*, 29(1):73–95, 2024.
- [54] Mohammadmehdi Ateai, Hyunmin Cheong, Daniele Grandi, Ye Wang, Nigel Morris, and Alexander Tessier. Elictron: An LLM agent-based simulation framework for design requirements elicitation. *CoRR*, abs/2404.16045, 2024.
- [55] Lezhi Ma, Shangqing Liu, Yi Li, Xiaofei Xie, and Lei Bu. Specgen: Automated generation of formal program specifications via large language models. *arXiv preprint arXiv:2401.08807*, 2024.
- [56] Chetan Arora, John Grundy, and Mohamed Abdelrazek. Advancing requirements engineering through generative ai: Assessing the role of llms. In *Generative AI for Effective Software Development*, pages 129–148. Springer, 2024.
- [57] Dongming Jin, Zhi Jin, Xiaohong Chen, and Chunhui Wang. MARE: multi-agents collaboration framework for requirements engineering. *CoRR*, abs/2405.03256, 2024.
- [58] Lilian Burdy, Yoonsik Cheon, David R Cok, Michael D Ernst, Joseph R Kiniry, Gary T Leavens, K Rustan M Leino, and Erik Poll. An overview of jml tools and applications. *International journal on software tools for technology transfer*, 7:212–232, 2005.
- [59] David R Cok. Openjml: Jml for java 7 by extending openjdk. In *NASA Formal Methods: Third International Symposium, NFM 2011, Pasadena, CA, USA, April 18-20, 2011. Proceedings 3*, pages 472–479. Springer, 2011.
- [60] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lema Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. Siren’s song in the AI ocean: A survey on hallucination in large language models. *CoRR*, abs/2309.01219, 2023.
- [61] Eric Zelikman, Qian Huang, Gabriel Poesia, Noah D. Goodman, and Nick Haber. Parsel: Algorithmic reasoning with language models by composing decompositions. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [62] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023.
- [63] Theo X. Olausson, Jeevana Priya Inala, Chenglong Wang, Jianfeng Gao, and Armando Solar-Lezama. Is self-repair a silver bullet for code generation?, 2024.

- [64] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryan W White, Doug Burger, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation, 2023.
- [65] Hanbin Wang, Zhenghao Liu, Shuo Wang, Ganqu Cui, Ning Ding, Zhiyuan Liu, and Ge Yu. Intervenor: Prompting the coding ability of large language models with the interactive chain of repair, 2024.
- [66] Noble Saji Mathews and Meiyappan Nagappan. Test-driven development for code generation. *CoRR*, abs/2402.13521, 2024.
- [67] Bin Lei, Yuchen Li, and Qiuwu Chen. Autocoder: Enhancing code large language model with aiev-instruct. *CoRR*, abs/2405.14906, 2024.
- [68] Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for “mind” exploration of large language model society, 2023.
- [69] Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, and Deheng Ye. More agents is all you need, 2024.
- [70] Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. Dynamic llm-agent network: An llm-agent collaboration framework with agent team optimization, 2023.
- [71] Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug, 2023.
- [72] Vadim Liventsev, Anastasiia Grishina, Aki Härmä, and Leon Moonen. Fully autonomous programming with large language models. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 1146–1155, 2023.
- [73] Zhao Tian, Junjie Chen, and Xiangyu Zhang. Test-case-driven programming understanding in large language models for better code generation, 2024.
- [74] Tal Ridnik, Dedy Kredo, and Itamar Friedman. Code generation with alphacodium: From prompt engineering to flow engineering, 2024.
- [75] Lily Zhong, Zilong Wang, and Jingbo Shang. LDB: A large language model debugger via verifying runtime execution step-by-step. *CoRR*, abs/2402.16906, 2024.
- [76] Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language agent tree search unifies reasoning acting and planning in language models. *CoRR*, abs/2310.04406, 2023.
- [77] Ajinkya Deshpande, Anmol Agarwal, Shashank Shet, Arun Iyer, Aditya Kanade, Ramakrishna Bairi, and Suresh Parthasarathy. Class-level code generation from natural language using iterative, tool-enhanced reasoning over repository, 2024.
- [78] Kechi Zhang, Huangzhao Zhang, Ge Li, Jia Li, Zhuo Li, and Zhi Jin. Toolcoder: Teach code generation models to use api search tools, 2023.
- [79] Shuyang Jiang, Yuhao Wang, and Yu Wang. Selfevolve: A code evolution framework via large language models, 2023.
- [80] Xiaoxue Ren, Xinyuan Ye, Dehai Zhao, Zhenchang Xing, and Xiaohu Yang. From misuse to mastery: Enhancing code generation with knowledge-driven ai chaining, 2023.
- [81] Yiheng Xu, Hongjin Su, Chen Xing, Boyu Mi, Qian Liu, Weijia Shi, Binyuan Hui, Fan Zhou, Yitao Liu, Tianbao Xie, Zhoujun Cheng, Siheng Zhao, Lingpeng Kong, Bailin Wang, Caiming Xiong, and Tao Yu. Lemur: Harmonizing natural language and code for language agents, 2023.
- [82] Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin. Codeagent: Enhancing code generation with tool-integrated agent systems for real-world repo-level coding challenges, 2024.
- [83] Sanyogita Piya and Allison Sullivan. Llm4td: Best practices for test driven development using large language models, 2023.
- [84] Dong Huang, Qingwen Bu, Yuhao Qing, and Heming Cui. Codcot: Tackling code syntax errors in cot reasoning for code generation, 2024.
- [85] Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. Executable code actions elicit better llm agents, 2024.
- [86] Xinyi He, Jiaru Zou, Yun Lin, Mengyu Zhou, Shi Han, Zejian Yuan, and Dongmei Zhang. CONLINE: complex code generation and refinement with online searching and correctness testing. *CoRR*, abs/2403.13583, 2024.
- [87] John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. Intercoder: Standardizing and benchmarking interactive coding with execution feedback, 2023.
- [88] Ramakrishna Bairi, Atharv Sonwane, Aditya Kanade, Vageesh D C, Arun Iyer, Suresh Parthasarathy, Sriram Rajamani, B. Ashok, and Shashank Shet. Codeplan: Repository-level coding using llms and planning, 2023.
- [89] Chong Wang, Jian Zhang, Yebo Feng, Tianlin Li, Weisong Sun, Yang Liu, and Xin Peng. Teaching code llms to use autocompletion tools in repository-level code generation, 2024.
- [90] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [91] Martin Josifoski, Lars Henning Klein, Maxime Peyrard, Yifei Li, Saibo Geng, Julian Paul Schnitzler, Yuxing Yao, Jiheng Wei, Debjit Paul, and Robert West. Flows: Building blocks of reasoning and collaborating AI. *CoRR*, abs/2308.01285, 2023.
- [92] Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. Mint: Evaluating llms in multi-turn interaction with tools and language feedback, 2024.
- [93] Hung Le, Hailin Chen, Amrita Saha, Akash Gokul, Doyen Sahoo, and Shafiq Joty. Codechain: Towards modular code generation through chain of self-revisions towards representative sub-modules, 2024.
- [94] Fangwen Mu, Lin Shi, Song Wang, Zhuohao Yu, Binquan Zhang, Chenxue Wang, Shichao Liu, and Qing Wang. Clarifygpt: Empowering llm-based code generation with intention clarification, 2023.
- [95] Dong Huang, Jie M. Zhang, Michael Luck, Qingwen Bu, Yuhao Qing, and Heming Cui. Agentcoder: Multi-agent-based code generation with iterative testing and optimisation, 2024.
- [96] Binfeng Xu, Xukun Liu, Hua Shen, Zeyu Han, Yuhao Li, Murong Yue, Zhiyuan Peng, Yuchen Liu, Ziyu Yao, and Dongkuan Xu. Gentopia: A collaborative platform for tool-augmented llms, 2023.
- [97] Yoichi Ishibashi and Yoshimasa Nishimura. Self-organized agents: A LLM multi-agent framework toward ultra large-scale code generation and optimization. *CoRR*, abs/2404.02183, 2024.
- [98] Md. Ashraf Islam, Mohammed Eunus Ali, and Md. Rizwan Parvez. Mapcoder: Multi-agent code generation for competitive problem solving. *CoRR*, abs/2405.11403, 2024.
- [99] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- [100] Zhiqiang Yuan, Junwei Liu, Qiancheng Zi, Mingwei Liu, Xin Peng, and Yiling Lou. Evaluating instruction-tuned large language models on code comprehension and generation. *arXiv preprint arXiv:2308.01240*, 2023.
- [101] Xueying Du, Geng Zheng, Kaixin Wang, Jiayi Feng, Wentai Deng, Mingwei Liu, Bihuan Chen, Xin Peng, Tao Ma, and Yiling Lou. Vul-rag: Enhancing llm-based vulnerability detection via knowledge-level rag. *arXiv preprint arXiv:2406.11147*, 2024.
- [102] Bhargavi Paranjape, Scott M. Lundberg, Sameer Singh, Hananeh Hajishirzi, Luke Zettlemoyer, and Marco Túlio Ribeiro. ART: automatic multi-step reasoning and tool-use for large language models. *CoRR*, abs/2303.09014, 2023.
- [103] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adria Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*, 2022.
- [104] Sihao Hu, Tiansheng Huang, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. Large language model-powered smart contract vulnerability detection: New perspectives. In *5th IEEE International Conference on Trust, Privacy and Security in Intelligent Systems and Applications, TPS-ISA 2023, Atlanta, GA, USA, November 1-4, 2023*, pages 297–306. IEEE, 2023.

- [105] Gang Fan, Xiaoheng Xie, Xunjin Zheng, Yinan Liang, and Peng Di. Static code analysis in the AI era: An in-depth exploration of the concept, function, and potential of intelligent code analysis agents. *CoRR*, abs/2310.08837, 2023.
- [106] Aida Radu and Sarah Nadi. A dataset of non-functional bugs. In Margaret-Anne D. Storey, Bram Adams, and Sonia Haiduc, editors, *Proceedings of the 16th International Conference on Mining Software Repositories, MSR 2019, 26-27 May 2019, Montreal, Canada*, pages 399–403. IEEE / ACM, 2019.
- [107] Yu Hao, Weiteng Chen, Ziqiao Zhou, and Weidong Cui. E&v: Prompting large language models to perform static analysis by pseudo-code execution and verification. *CoRR*, abs/2312.08477, 2023.
- [108] *Clang*, 2024. <https://clang.llvm.org/>.
- [109] *syzbot*, 2024. <https://syzkaller.appspot.com/upstream>.
- [110] Yuqiang Sun, Daoyuan Wu, Yue Xue, Han Liu, Wei Ma, Lyuye Zhang, Miaolei Shi, and Yang Liu. Llm4vuln: A unified evaluation framework for decoupling and enhancing llms' vulnerability reasoning. *CoRR*, abs/2401.16185, 2024.
- [111] Zhenyu Mao, Jialong Li, Munan Li, and Kenji Tei. Multi-role consensus through llms discussions for vulnerability detection. *CoRR*, abs/2403.14274, 2024.
- [112] Zhen Li, Deqing Zou, Shouhuai Xu, Hai Jin, Yawei Zhu, and Zhaoxuan Chen. Sysevr: A framework for using deep learning to detect software vulnerabilities. *IEEE Trans. Dependable Secur. Comput.*, 19(4):2244–2258, 2022.
- [113] Ziyang Li, Saikat Dutta, and Mayur Naik. Llm-assisted static analysis for detecting security vulnerabilities, 2024.
- [114] Pavel Avgustinov, Oege De Moor, Michael Peyton Jones, and Max Schäfer. Ql: Object-oriented queries on relational data. In *30th European Conference on Object-Oriented Programming (ECOOP 2016)*. Schloss-Dagstuhl-Leibniz Zentrum für Informatik, 2016.
- [115] Haonan Li, Yu Hao, Yizhuo Zhai, and Zhiyun Qian. Enhancing static analysis for practical bug detection: An llm-integrated approach. *Proc. ACM Program. Lang.*, 8(OOPSLA1):474–499, 2024.
- [116] Yizhuo Zhai, Yu Hao, Hang Zhang, Daimeng Wang, Chengyu Song, Zhiyun Qian, Mohsen Lesani, Srikanth V Krishnamurthy, and Paul Yu. Ubitec: a precise and scalable method to detect use-before-initialization bugs in linux kernel. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 221–232, 2020.
- [117] Taghi M Khoshgoftaar and Naeem Seliya. Comparative assessment of software quality classification techniques: An empirical case study. *Empirical Software Engineering*, 9:229–257, 2004.
- [118] Lingwei Li, Li Yang, Huaxi Jiang, Jun Yan, Tiejian Luo, Zihan Hua, Geng Liang, and Chun Zuo. Auger: automatically generating review comments with pre-training models. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 1009–1021, 2022.
- [119] Daniel Tang, Zhenghan Chen, Kisub Kim, Yewei Song, Haoye Tian, Saad Ezzini, Yongfeng Huang, Jacques Klein, and Tegawendé F. Bissyandé. Codeagent: Collaborative agents for software engineering. *CoRR*, abs/2402.02172, 2024.
- [120] Zeeshan Rasheed, Malik Abdul Sami, Muhammad Waseem, Kai-Kristian Kemell, Xiaofeng Wang, Anh Nguyen, Kari Systä, and Pekka Abrahamsson. Ai-powered code review with llms: Early results. *CoRR*, abs/2404.18496, 2024.
- [121] Nalin Wadhwa, Jui Pradhan, Atharv Sonwane, Surya Prakash Sahu, Nagarajan Natarajan, Aditya Kanade, Suresh Parthasarathy, and Sriram K. Rajamani. Frustrated with code quality issues? llms can help! *CoRR*, abs/2309.12938, 2023.
- [122] Zhiqiang Yuan, Yiling Lou, Mingwei Liu, Shiji Ding, Kaixin Wang, Yixuan Chen, and Xin Peng. No more manual tests? evaluating and improving chatgpt for unit test generation, 2024.
- [123] Max Schäfer, Sarah Nadi, Aryaz Eghbali, and Frank Tip. An empirical evaluation of using large language models for automated unit test generation, 2023.
- [124] Zhuokui Xie, Yinghao Chen, Chen Zhi, Shuiguang Deng, and Jianwei Yin. Chatunitest: a chatgpt-based automated unit test generation tool. *CoRR*, abs/2305.04764, 2023.
- [125] Chen Yang, Junjie Chen, Bin Lin, Jianyi Zhou, and Ziqi Wang. Enhancing llm-based test generation for hard-to-cover branches via program analysis, 2024.
- [126] Juan Altmayer Pizzorno and E. Berger. Coverup: Coverage-guided llm-based test generation. *ArXiv*, abs/2403.16218, 2024.
- [127] Arghavan Moradi Dakhel, Amin Nikanjam, Vahid Majdinasab, Foutse Khomh, and Michel C. Desmarais. Effective test generation using pre-trained large language models and mutation testing, 2023.
- [128] Chunqiu Steven Xia and Lingming Zhang. Conversational automated program repair. *arXiv preprint arXiv:2301.13246*, 2023.
- [129] Juan Altmayer Pizzorno and Emery D. Berger. Slipcover: Near zero-overhead code coverage for python. In *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA '23*. ACM, July 2023.
- [130] Chenyuan Yang, Zijie Zhao, and Lingming Zhang. Kernelgpt: Enhanced kernel fuzzing via large language models, 2023.
- [131] *Syzkaller*. <https://github.com/google/syzkaller/>.
- [132] *The LLVM Compiler Infrastructure*. <https://llvm.org>.
- [133] Chenyuan Yang, Yinlin Deng, Runyu Lu, Jiayi Yao, Jiawei Liu, Reyhaneh Jabbarvand, and Lingming Zhang. White-box compiler fuzzing empowered by large language models, 2023.
- [134] Haoxin Tu, Zhide Zhou, He Jiang, Imam Nur Bani Yusuf, Yuxian Li, and Lingxiao Jiang. Isolating compiler bugs by generating effective witness programs with large language models. *IEEE Transactions on Software Engineering*, page 1–20, 2024.
- [135] *OCLint*. <https://github.com/oclint/>.
- [136] Christian D Newman, Tessandra Sage, Michael L Collard, Hakam W Alomari, and Jonathan I Maletic. srslslic: A tool for efficient static forward slicing. In *Proceedings of the 38th International Conference on Software Engineering Companion*, pages 621–624, 2016.
- [137] *Gcov*, 2023. <https://gcc.gnu.org/onlinedocs/gcc/Gcov.html>.
- [138] *Frama-C*. <https://www.frama-c.com/>.
- [139] Zhe Liu, Chunyang Chen, Junjie Wang, Mengzhuo Chen, Boyu Wu, Xing Che, Dandan Wang, and Qing Wang. Make llm a testing expert: Bringing human-like interaction to mobile gui testing via functionality-aware decisions. *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, 2023.
- [140] *virtualbox*, 2023. <https://www.virtualbox.org/>.
- [141] *pyvbox*, 2023. <https://pypi.org/project/pyvbox/>.
- [142] *Python wrapper of Android uiautomator test tool*, 2021. <https://github.com/xiaocong/uiautomator>.
- [143] *Android Debug Bridge (adb) - Android Developers*, 2023. <https://developer.android.com/studio/command-line/adb/>.
- [144] Juyeon Yoon, Robert Feldt, and Shin Yoo. Autonomous large language model agents enabling intent-driven mobile gui testing. *ArXiv*, abs/2311.08649, 2023.
- [145] Maryam Taeb, Amanda Swearngin, Eldon Schoop, Ruijia Cheng, Yue Jiang, and Jeffrey Nichols. Axnav: Replaying accessibility tests from natural language. In *Proceedings of the CHI Conference on Human Factors in Computing Systems, CHI '24*. ACM, May 2024.
- [146] *Genymotion – Android Emulator for app testing*, 2023. <https://www.genymotion.com/>.
- [147] *Android UiAutomator2 Python Wrapper*, 2023. <https://github.com/openatx/uiautomator2/>.
- [148] Zhitao Wang, Wei Wang, Zirao Li, Long Wang, Can Yi, Xinjie Xu, Luyang Cao, Hanjing Su, Shouzhi Chen, and Jun Zhou. Xuat-copilot: Multi-agent collaborative system for automated user acceptance testing with large language model, 2024.
- [149] Alix Decrop, Gilles Perrouin, Mike Papadakis, Xavier Devroey, and Pierre-Yves Schobbens. You can rest now: Automated specification inference and black-box testing of restful apis with large language models, 2024.
- [150] Chunqiu Steven Xia, Matteo Paltenghi, Jia Le Tian, Michael Pradel, and Lingming Zhang. Fuzz4all: Universal fuzzing with large language models, 2024.
- [151] Gelei Deng, Yi Liu, Víctor Mayoral-Vilches, Peng Liu, Yuekang Li, Yuan Xu, Tianwei Zhang, Yang Liu, Martin Pinzger, and Stefan Rass. Pentestgpt: An llm-empowered automatic penetration testing tool, 2024.
- [152] *Metasploit framework*. <https://www.metasploit.com/>.
- [153] Richard Fang, Rohan Bindu, Akul Gupta, and Daniel Kang. Llm agents can autonomously exploit one-day vulnerabilities, 2024.
- [154] Zhe Liu, Chunyang Chen, Junjie Wang, Xing Che, Yuekai Huang, Jun Hu, and Qing Wang. Fill in the blank: Context-aware automated text input generation for mobile GUI testing. In *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*, pages 1355–1367. IEEE, 2023.
- [155] Hao Wen, Hongming Wang, Jiaxuan Liu, and Yuanchun Li. Droidbot-gpt: Gpt-powered UI automation for android. *CoRR*, abs/2304.07061, 2023.

- [156] Sjouke Mauw and Martijn Oostdijk. Foundations of attack trees. In *Information Security and Cryptology-ICISC 2005: 8th International Conference, Seoul, Korea, December 1-2, 2005, Revised Selected Papers 8*, pages 186–198. Springer, 2006.
- [157] W Eric Wong, Ruizhi Gao, Yihao Li, Rui Abreu, and Franz Wotawa. A survey on software fault localization. *IEEE Transactions on Software Engineering*, 42(8):707–740, 2016.
- [158] Luca Gazzola, Daniela Micucci, and Leonardo Mariani. Automatic software repair: A survey. In *Proceedings of the 40th International Conference on Software Engineering*, pages 1219–1219, 2018.
- [159] Cheryl Lee, Chunqiu Steven Xia, Jen-tse Huang, Zhouxiu Xing, Lingming Zhang, and Michael R. Lyu. A unified debugging approach via llm-based multi-agent synergy. *CoRR*, abs/2404.17153, 2024.
- [160] Yihao Qin, Shangwen Wang, Yiling Lou, Jinhao Dong, Kaixin Wang, Xiaoling Li, and Xiaoguang Mao. Agentfl: Scaling llm-based fault localization to project-level context. *CoRR*, abs/2403.16362, 2024.
- [161] *tree-sitter*. <https://github.com/tree-sitter/tree-sitter/>.
- [162] Zefan Wang, Zichuan Liu, Yingying Zhang, Aoxiao Zhong, Lunting Fan, Lingfei Wu, and Qingsong Wen. Reagent: Cloud root cause analysis by autonomous agents with tool-augmented large language models. *CoRR*, abs/2310.16340, 2023.
- [163] Sungmin Kang, Gabin An, and Shin Yoo. A preliminary evaluation of llm-based fault localization. *CoRR*, abs/2308.05487, 2023.
- [164] Xia Li, Wei Li, Yuqun Zhang, and Lingming Zhang. Deepfl: Integrating multiple fault diagnosis dimensions for deep fault localization. In *Proceedings of the 28th ACM SIGSOFT international symposium on software testing and analysis*, pages 169–180, 2019.
- [165] Zimin Chen, Steve Kommrusch, Michele Tufano, Louis-Noël Pouchet, Denys Poshyvanyk, and Martin Monperrus. Sequencer: Sequence-to-sequence learning for end-to-end program repair. *IEEE Transactions on Software Engineering*, 47(9):1943–1959, 2019.
- [166] Chunqiu Steven Xia and Lingming Zhang. Automated program repair via conversation: Fixing 162 out of 337 bugs for \$0.42 each using chatgpt. In *ISSTA*, 2024.
- [167] René Just, Darioush Jalali, and Michael D Ernst. Defects4j: A database of existing faults to enable controlled testing studies for java programs. In *Proceedings of the 2014 international symposium on software testing and analysis*, pages 437–440, 2014.
- [168] Derrick Lin, James Koppel, Angela Chen, and Armando Solar-Lezama. Quixbugs: a multi-lingual program repair benchmark set based on the quixey challenge. In Gail C. Murphy, editor, *Proceedings Companion of the 2017 ACM SIGPLAN International Conference on Systems, Programming, Languages, and Applications: Software for Humanity, SPLASH 2017, Vancouver, BC, Canada, October 23 - 27, 2017*, pages 55–56. ACM, 2017.
- [169] Dávid Hidvégi, Khashayar Etemadi, Sofia Bobadilla, and Martin Monperrus. Cigar: Cost-efficient program repair with llms. *CoRR*, abs/2402.06598, 2024.
- [170] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021.
- [171] Islem Bouzenia, Premkumar T. Devanbu, and Michael Pradel. Repairagent: An autonomous, llm-based agent for program repair. *CoRR*, abs/2403.17134, 2024.
- [172] Sungmin Kang, Bei Chen, Shin Yoo, and Jian-Guang Lou. Explainable automated debugging via large language model-driven scientific debugging. *CoRR*, abs/2304.02195, 2023.
- [173] Ratnadira Widayarsi, Sheng Qin Sim, Camellia Lok, Haodi Qi, Jack Phan, Qijin Tay, Constance Tan, Fiona Wee, Jodie Ethelda Tan, Yuheng Yieh, et al. Bugsinpy: a database of existing bugs in python programs to enable controlled testing and debugging studies. In *Proceedings of the 28th ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering*, pages 1556–1560, 2020.
- [174] Lyuye Zhang, Kaixuan Li, Kairan Sun, Daoyuan Wu, Ye Liu, Haoye Tian, and Yang Liu. ACFIX: guiding llms with mined common RBAC practices for context-aware repair of access control vulnerabilities in smart contracts. *CoRR*, abs/2403.06838, 2024.
- [175] Yang Chen. Flakiness repair in the era of large language models. In *Proceedings of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion Proceedings*, pages 441–443, 2024.
- [176] *International Dataset of Flaky tests*. <https://github.com/TestingResearchIllinois/idoft>.
- [177] Peilun Zhang, Yanjie Jiang, Anjiang Wei, Victoria Stodden, Darko Marinov, and August Shi. Domain-specific fixes for flaky tests with wrong assumptions on underdetermined specifications. In *43rd IEEE/ACM International Conference on Software Engineering, ICSE 2021, Madrid, Spain, 22-30 May 2021*, pages 50–61. IEEE, 2021.
- [178] Chengpeng Li, Chenguang Zhu, Wenxi Wang, and August Shi. Repairing order-dependent flaky tests via test generation. In *44th IEEE/ACM International Conference on Software Engineering, ICSE 2022, Pittsburgh, PA, USA, May 25-27, 2022*, pages 1881–1892. ACM, 2022.
- [179] Andreas Zeller. *Why Programs Fail - A Guide to Systematic Debugging, 2nd Edition*. Academic Press, 2009.
- [180] Yiling Lou, Ali Ghanbari, Xia Li, Lingming Zhang, Haotian Zhang, Dan Hao, and Lu Zhang. Can automated program repair refine fault localization? a unified debugging approach. In Sarfraz Khurshid and Corina S. Pasareanu, editors, *ISSTA '20: 29th ACM SIGSOFT International Symposium on Software Testing and Analysis, Virtual Event, USA, July 18-22, 2020*, pages 75–87. ACM, 2020.
- [181] Samuel Benton, Xia Li, Yiling Lou, and Lingming Zhang. Evaluating and improving unified debugging. *IEEE Trans. Software Eng.*, 48(11):4692–4716, 2022.
- [182] Samuel Benton, Xia Li, Yiling Lou, and Lingming Zhang. On the effectiveness of unified debugging: An extensive study on 16 program repair systems. In *35th IEEE/ACM International Conference on Automated Software Engineering, ASE 2020, Melbourne, Australia, September 21-25, 2020*, pages 907–918. IEEE, 2020.
- [183] Yuzhe Cai, Shaoguang Mao, Wenshan Wu, Zehua Wang, Yaobo Liang, Tao Ge, Chenfei Wu, Wang You, Ting Song, Yan Xia, Jonathan Tien, and Nan Duan. Low-code LLM: visual programming over llms. *CoRR*, abs/2304.08103, 2023.
- [184] Zhenchang Xing, Qing Huang, Yu Cheng, Liming Zhu, Qinghua Lu, and Xiwei Xu. Prompt sapper: Llm-empowered software engineering infrastructure for ai-native services. *CoRR*, abs/2306.02230, 2023.
- [185] Yashar Talebirad and Amirhossein Nadiri. Multi-agent collaboration: Harnessing the power of intelligent LLM agents. *CoRR*, abs/2306.03314, 2023.
- [186] Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. Communicative agents for software development. *CoRR*, abs/2307.07924, 2023.
- [187] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin Wang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. Metagpt: Meta programming for a multi-agent collaborative framework, 2023.
- [188] Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chen Qian, Chi-Min Chan, Yujia Qin, Yaxi Lu, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors in agents. *CoRR*, abs/2308.10848, 2023.
- [189] Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Jaward Sesay, Börje F. Karlsson, Jie Fu, and Yemin Shi. Autoagents: A framework for automatic agent generation. *CoRR*, abs/2309.17288, 2023.
- [190] Chen Qian, Yufan Dang, Jiahao Li, Wei Liu, Weize Chen, Cheng Yang, Zhiyuan Liu, and Maosong Sun. Experiential co-learning of software-developing agents. *CoRR*, abs/2312.17025, 2023.
- [191] Simiao Zhang, Jiaping Wang, Guoliang Dong, Jun Sun, Yueling Zhang, and Geguang Pu. Experimenting a new programming practice with llms. *CoRR*, abs/2401.01062, 2024.
- [192] Mohamad Fakhri, Rahul Dharmaji, Yasamin Moghaddas, Gustavo Quiros, Oluwatosin Ogundare, and Mohammad Abdullah Al Faruque. LLM4PLC: harnessing large language models for verifiable programming of plcs in industrial control systems. In *Pro-*

- ceedings of the 46th International Conference on Software Engineering: Software Engineering in Practice, ICSE-SEIP 2024, Lisbon, Portugal, April 14-20, 2024*, pages 192–203. ACM, 2024.
- [193] Zeeshan Rasheed, Muhammad Waseem, Mika Saari, Kari Systä, and Pekka Abrahamsson. Codepori: Large scale model for autonomous software development by using multi-agents. *CoRR*, abs/2402.01411, 2024.
- [194] Feng Lin, Dong Jae Kim, and Tse-Hsun Chen. When llm-based code generation meets the software development process. *CoRR*, abs/2403.15852, 2024.
- [195] Daoguang Zan, Ailun Yu, Wei Liu, Dong Chen, Bo Shen, Wei Li, Yafen Yao, Yongshun Gong, Xiaolin Chen, Bei Guan, Zhiguang Yang, Yongji Wang, Qianxiang Wang, and Lizhen Cui. Codes: Natural language to code repository via multi-layer sketch. *CoRR*, abs/2403.16443, 2024.
- [196] Chen Qian, Jiahao Li, Yufan Dang, Wei Liu, Yifei Wang, Zihao Xie, Weize Chen, Cheng Yang, Yingli Zhang, Zhiyuan Liu, and Maosong Sun. Iterative experience refinement of software-developing agents. *CoRR*, abs/2405.04219, 2024.
- [197] Zhuoyun Du, Chen Qian, Wei Liu, Zihao Xie, Yifei Wang, Yufan Dang, Weize Chen, and Cheng Yang. Multi-agent software development through cross-team collaboration, 2024.
- [198] Minh Huynh Nguyen, Thang Phan Chau, Phong X Nguyen, and Nghi DQ Bui. Agilecoder: Dynamic collaborative agents for software development based on agile methodology. *arXiv preprint arXiv:2406.11912*, 2024.
- [199] W. W. Royce. Managing the development of large software systems: Concepts and techniques. In William E. Riddle, Robert M. Balzer, and Kouichi Kishida, editors, *Proceedings, 9th International Conference on Software Engineering, Monterey, California, USA, March 30 - April 2, 1987*, pages 328–339. ACM Press, 1987.
- [200] HA Thakur and SD Maurya. A survey on incremental software development life cycle model. *Int. J. Eng. Technol. Comput. Res*, 3(2):102–106, 2016.
- [201] Ivar Jacobson, Grady Booch, and James Rumbaugh. The unified process. *Ieee Software*, 16(3):96, 1999.
- [202] Mikio Aoyama. Agile software process model. In *Proceedings Twenty-First Annual International Computer Software and Applications Conference (COMPSAC'97)*, pages 454–459. IEEE, 1997.
- [203] Yihong Dong, Jiazhen Ding, Xue Jiang, Ge Li, Zhuo Li, and Zhi Jin. Codescore: Evaluating code generation by learning code execution, 2023.
- [204] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models, 2021.
- [205] Wei Tao, Yucheng Zhou, Wenqiang Zhang, and Yu Cheng. MAGIS: llm-based multi-agent framework for github issue resolution. *CoRR*, abs/2403.17927, 2024.
- [206] Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. Autocoderover: Autonomous program improvement. *CoRR*, abs/2404.05427, 2024.
- [207] John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering. *arXiv preprint arXiv:2405.15793*, 2024.
- [208] Dong Chen, Shaoxin Lin, Muhan Zeng, Daoguang Zan, Jiang-Gang Wang, Anton Cheshkov, Jun Sun, Hao Yu, Guoliang Dong, Artem Aliev, et al. Coder: Issue resolving with multi-agent and task graphs. *arXiv preprint arXiv:2406.01304*, 2024.
- [209] Yingwei Ma, Qingping Yang, Rongyu Cao, Binhua Li, Fei Huang, and Yongbin Li. How to understand whole software repository? *arXiv preprint arXiv:2406.01422*, 2024.
- [210] Daman Arora, Atharv Sonwane, Nalin Wadhwa, Abhav Mehrotra, Saiteja Utpala, Ramakrishna Bairi, Aditya Kanade, and Nagarajan Natarajan. Masai: Modular architecture for software-engineering ai agents, 2024.
- [211] Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demystifying llm-based software engineering agents. *CoRR*, abs/2407.01489, 2024.
- [212] Stephen E Robertson and Steve Walker. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In *SIGIR'94: Proceedings of the Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, organised by Dublin City University*, pages 232–241. Springer, 1994.
- [213] Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *CoRR*, abs/2310.06770, 2023.
- [214] SWE-bench Lite, 2024. <https://www.swebench.com/lite.html>.
- [215] *Introducing SWE-bench Verified*, 2024. <https://openai.com/index/introducing-swe-bench-verified/>.
- [216] *Function Calling and other API updates*, 2023. <https://openai.com/index/function-calling-and-other-api-updates/>.
- [217] GPT-3.5, 2023. <https://platform.openai.com/docs/models/gpt-3-5-turbo>.
- [218] GPT-4, 2023. <https://platform.openai.com/docs/models/gpt-4-and-gpt-4-turbo>.
- [219] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- [220] Yuzhe Cai, Shaoguang Mao, Wenshan Wu, Zehua Wang, Yaobo Liang, Tao Ge, Chenfei Wu, Wang You, Ting Song, Yan Xia, et al. Low-code llm: Visual programming over llms. *arXiv preprint arXiv:2304.08103*, 2, 2023.
- [221] Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 5484–5495. Association for Computational Linguistics, 2021.
- [222] Iain D Craig. Blackboard systems. *Artificial Intelligence Review*, 2(2):103–118, 1988.
- [223] Hatice Koç, Ali Mert Erdoğan, Yousef Barjakly, and Serhat Peker. Uml diagrams in software engineering research: a systematic literature review. In *Proceedings*, volume 74, page 13. MDPI, 2021.
- [224] Jun Tang, Zhibo Yang, Yongpan Wang, Qi Zheng, Yongchao Xu, and Xiang Bai. Seglink++: Detecting dense and arbitrary-shaped scene text by instance-aware component grouping. *Pattern recognition*, 96:106954, 2019.
- [225] Zhuang Liu, Hanzhi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022.
- [226] Cheng Wen, Jialun Cao, Jie Su, Zhiwu Xu, Shengchao Qin, Mengda He, Haokun Li, Shing-Chi Cheung, and Cong Tian. Enchanting program specification synthesis by large language models using static analysis and program verification, 2024.
- [227] Qinyu Luo, Yining Ye, Shihao Liang, Zhong Zhang, Yujia Qin, Yaxi Lu, Yesai Wu, Xin Cong, Yankai Lin, Yingli Zhang, Xiaoyin Che, Zhiyuan Liu, and Maosong Sun. Repoagent: An llm-powered open-source framework for repository-level code documentation generation, 2024.
- [228] *DuckDuckGo*. <https://duckduckgo.com/>.
- [229] *SerpApi*. <https://serpapi.com/>.
- [230] Xinyi He, Jiuru Zou, Yun Lin, Mengyu Zhou, Shi Han, Zejian Yuan, and Dongmei Zhang. Cocost: Automatic complex code generation with online searching and correctness testing, 2024.
- [231] Jun Tang, Zhibo Yang, Yongpan Wang, Qi Zheng, Yongchao Xu, and Xiang Bai. Seglink++: Detecting dense and arbitrary-shaped scene text by instance-aware component grouping. *Pattern Recognition*, 96:106954, 2019.
- [232] Xiaoyi Zhang, Lilian De Greef, Amanda Swearngin, Samuel White, Kyle Murray, Lisa Yu, Qi Shan, Jeffrey Nichols, Jason Wu, Chris Fleizach, et al. Screen recognition: Creating accessibility metadata for mobile applications from pixels. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–15, 2021.
- [233] Zhuang Liu, Hanzhi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022.
- [234] T. J. Parr and R. W. Quong. Antlr: A predicated-ll(k) parser generator. *Software: Practice and Experience*, 25(7):789–810, 1995.
- [235] *Jedi*. <https://github.com/davidhalter/jedi/>.
- [236] *EclipseJDTLS*. <https://github.com/eclipse-jdtls/eclipse-jdt.ls>.
- [237] *Black*. <https://github.com/psf/black>.
- [238] Roberto Cavada, Alessandro Cimatti, Michele Dorigatti, Alberto Griggio, Alessandro Mariotti, Andrea Micheli, Sergio Mover, Marco Roveri, and Stefano Tonetta. The nuxmv symbolic model checker. In Armin Biere and Roderick Bloem, editors, *Computer*

- Aided Verification*, pages 334–342, Cham, 2014. Springer International Publishing.
- [239] Josselin Feist, Gustavo Grieco, and Alex Groce. Slither: a static analysis framework for smart contracts. In *2019 IEEE/ACM 2nd International Workshop on Emerging Trends in Software Engineering for Blockchain (WETSEB)*, pages 8–15. IEEE, 2019.
- [240] *GPT-4*, 2024. <https://docs.oracle.com/javase/8/docs/api/java/lang/instrument/package-summary.html>.
- [241] Stephan Lukasczyk and Gordon Fraser. Pynguin: automated unit test generation for python. In *Proceedings of the ACM/IEEE 44th International Conference on Software Engineering: Companion Proceedings, ICSE '22*. ACM, May 2022.
- [242] Konrad Hałas. Mutpy: a mutation testing tool for python 3.x source code, 2019. <https://github.com/mutpy/mutpy>.
- [243] José Campos, André Ribeiro, Alexandre Perez, and Rui Abreu. Gzoltar: an eclipse plug-in for testing and debugging. In *Proceedings of the 27th IEEE/ACM international conference on automated software engineering*, pages 378–381, 2012.
- [244] *Git*. <https://git-scm.com/>.
- [245] Ali Nouri, Beatriz Cabrero Daniel, Fredrik Törner, Håkan Siven-crona, and Christian Berger. Engineering safety requirements for autonomous driving with large language models. *CoRR*, abs/2403.16289, 2024.
- [246] *OpenAI: Introducing ChatGPT*, 2022. <https://openai.com/blog/chatgpt/>.
- [247] 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. Deepseek-coder: When the large language model meets programming - the rise of code intelligence. *CoRR*, abs/2401.14196, 2024.
- [248] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cas-sano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. Starcoder 2 and the stack v2: The next generation. *arXiv preprint arXiv:2402.19173*, 2024.