# Designing LLM-based Multi-Agent Systems for Software Engineering Tasks: Quality Attributes, Design Patterns and Rationale

YANGXIAO CAI, School of Computer Science, Wuhan University, China
RUIYIN LI, School of Computer Science, Wuhan University, China
PENG LIANG, School of Computer Science, Wuhan University, China
MOJTABA SHAHIN, School of Computing Technologies, RMIT University, Australia
ZENGYANG LI, School of Computer Science, Central China Normal University, China

As the complexity of Software Engineering (SE) tasks continues to escalate, Multi-Agent Systems (MASs) have emerged as a focal point of research and practice due to their autonomy and scalability. Furthermore, through leveraging the reasoning and planning capabilities of Large Language Models (LLMs), the application of LLM-based MASs in the field of SE is garnering increasing attention. However, there is no dedicated study that systematically explores the design of LLM-based MASs, including the Quality Attributes (QAs) on which the designers mainly focus, the design patterns used by the designers, and the rationale guiding the design of LLM-based MASs for SE tasks. To this end, we conducted a study to identify the QAs that LLM-based MASs for SE tasks focus on, the design patterns used in the MASs, and the design rationale for the MASs. We collected 94 papers on LLM-based MASs for SE tasks as the source. Our study shows that: (1) *Code Generation* is the most common SE task solved by LLM-based MASs among ten identified SE tasks, (2) *Functional Suitability* is the QA on which designers of LLM-based MASs pay the most attention, (3) *Role-Based Cooperation* is the design pattern most frequently employed among 16 patterns used to construct LLM-based MASs, and (4) *Improving the Quality of Generated Code* is the most common rationale behind the design of LLM-based MASs. Based on the study results, we presented the implications for the design of LLM-based MASs to support SE tasks.

CCS Concepts: • **Software and its engineering** → **Designing software**.

Additional Key Words and Phrases: Multi-Agent, Software Design, Large Language Model, Software Engineering Task

## 1 Introduction

While traditional Multi-Agent Systems (MASs) face limitations in handling complex decision-making tasks [8], Large Language Models (LLMs) can enhance the reasoning [12] and planning

Authors' Contact Information: Yangxiao Cai, yangxiaocai@whu.edu.cn, School of Computer Science, Wuhan University, Wuhan, China; Ruiyin Li, ryli_cs@whu.edu.cn, School of Computer Science, Wuhan University, Wuhan, China; Peng Liang, liangp@whu.edu.cn, School of Computer Science, Wuhan University, Wuhan, China; Mojtaba Shahin, School of Computing Technologies, RMIT University, Australia, mojtaba.shahin@rmit.edu.au; Zengyang Li, School of Computer Science, Central China Normal University, China, zengyangli@ccnu.edu.cn.

capabilities of MASs [23]. LLM-based MASs consist of multiple autonomous agents that collaborate through communication and responsibility specialization to tackle complex tasks given by users and simulate problem-solving environments [9]. LLM-based MASs have demonstrated significant potential in addressing Software Engineering (SE) tasks and have become a focal point of research and practice in SE [10]. Many researchers have designed LLM-based MASs to address specific SE tasks [10]. For example, Hong *et al.* [11] proposed an LLM-based MAS named MetaGPT, which is capable of automatically developing an entire software system. Jin *et al.* [14] proposed an LLM-based MAS named MARE (Multi-Agents Collaboration Framework for Requirements Engineering), which could complete the process of requirements engineering.

Effective design is paramount to the success of LLM-based MASs, as it directly shapes key quality attributes such as reliability, maintainability, scalability, and safety, and profoundly influences collaborative behaviors of agents (e.g., volunteering, conformity, and destructive actions), which in turn affect overall task-solving efficiency among agents in MASs [6]. Therefore, a principled understanding of design rationale is essential not only for constructing high-quality MASs but also for navigating trade-offs among competing quality attributes (e.g., correctness vs. efficiency) [15]. Despite the increasing interest and adoption of LLM-based MASs for various SE tasks, no existing study has systematically investigated the specific Quality Attributes (QAs), design patterns, and underlying design rationale that inform their construction. To address this gap, our **goal** is to identify QAs and design patterns that practitioners prioritize when building LLM-based MASs, thereby highlighting best practices and informing future system designs. Moreover, explicitly articulating the design rationale is critical for transparency of design decision-making and design knowledge transfer.

In this study, we empirically investigate quality-driven design choices in SE-focused LLM-based MASs. Our findings offer actionable guidance for practitioners and contribute a deeper, evidence-based understanding of how such systems are designed to be more robust, reliable, and effective. Our **study results** show that: (1) *Code Generation* is the most common SE task solved by LLM-based MASs among ten identified SE tasks, (2) *Functional Suitability* is the QA most frequently prioritized by the designers of MASs, (3) *Role-Based Cooperation* is the most commonly adopted design pattern among all 16 patterns we identified for MASs construction, and (4) *Improving the Quality of Generated Code* is the predominant rationale behind the design of LLM-based MASs. The analysis of QAs, design patterns, and design rationale was conducted specifically in the context of how LLM-based MASs are designed to complete these SE tasks, which ensures a more coherent logical flow from problem context to architectural analysis. The **main contributions** of this work:

- We collected a comprehensive set of SE tasks solved by LLM-based MASs. Moreover, we identified the key QAs that received the most attention from designers.
- We identified design patterns commonly used by designers for constructing LLM-based MASs to address specific SE tasks, and extracted underlying design rationale guiding their development.
- We established mapping relationships among SE tasks addressed by LLM-based MASs, QAs considered by designers, design patterns utilized in constructing LLM-based MASs for SE tasks, along with the design rationale.

The rest of this paper is structured as follows: Section 2 reviews the related work. Section 3 presents the Research Questions (RQs) and the research process conducted in this study. Section 4 presents the results of this study. Section 5 interprets the study results and discusses their implications. Section 6 outlines the potential threats to validity, and Section 7 concludes this work with future research directions.

## 2  Related Work

In this section, we first review LLM-based MASs specially designed to address SE tasks (Section 2.1), and introduce the studies that explore the characteristics of these MASs (Section 2.2). Then, we present the recent literature surveys that focus on LLM-based MASs (Section 2.3). We compare recent studies with our work to highlight the research gap (Section 2.4).

### 2.1  LLM-based MASs for SE Tasks

As LLMs continue to evolve at a rapid pace, SE researchers and practitioners are increasingly eager to employ the reasoning and decision-making abilities of LLMs to support SE tasks. To better exploit the potential of LLMs in complex, multi-step SE workflows (particularly those requiring coordination), LLM-based MASs have been designed to address SE tasks, such as requirements engineering, code generation, and fault localization. Arora *et al.* [1] developed an LLM-based MAS aimed at enhancing the efficiency and accuracy of requirements engineering tasks. This MAS reveals the significant potential of applying LLMs to requirements elicitation, analysis, specification, and validation. Li *et al.* [16] proposed *MAAD*, a knowledge-driven MAS specifically built for automated architecture design. *MAAD* assigns the architecture design task to four role-dedicated agents with external knowledge bases, and the four agents collaboratively generate architecture design by decomposing a given Software Requirements Specifications (SRS) into corresponding artifacts. Dong *et al.* [7] designed an LLM-based MAS based on a self-collaboration framework to achieve high-quality code generation. Their framework enables the MAS to solve complex repository-level tasks that are not readily solved by a single LLM agent. Pan *et al.* [19] proposed a novel optimization pipeline for LLM-based MASs based on self-generated textual feedback, named *CodeCoR*, which has a multi-phase workflow and iterative code repair progress. Its self-reflective mechanism can evaluate and refine the outputs of each agent. Some LLM-based MASs are also designed to complete specific SE tasks, such as code translation and software deployment. Yuan *et al.* [29] designed an LLM-based MAS named *TRANSAGENT* to improve code translation accuracy through the collaboration of four specialized agents, which could be used in software migration, system refactoring, and cross-platform development. By aligning the execution behavior of the source and target programs, *TRANSAGENT* effectively localizes faulty code blocks, thereby narrowing the error-fixing scope and reducing debugging complexity. Besides MASs tailored for specific SE tasks, an increasing number of studies have constructed LLM-based MASs to support holistic SE tasks such as end-to-end development and end-to-end maintenance. Qian *et al.* [20] developed an LLM-based MAS named *ChatDev*, which assigns multiple agents different roles in the software development lifecycle, such as design, coding, and testing to develop a complete software system.

### 2.2  Characteristics on LLM-based MASs

With the gradual adoption of LLM-based MASs in various SE tasks, many researchers have conducted empirical studies to examine their characteristics, including system architectures, inter-agent coordination mechanisms, and evaluation protocols. Shen *et al.* [22] designed a two-stage textual feedback optimization pipeline for LLM-based MASs focusing on software development. In the first stage, the MAS leverages its own failure explanations to identify underperforming agents. In the second stage, targeted prompt adjustments are applied to individual agents based on these explanations. The two-stage textual feedback pipeline provides a practical architecture pattern for runtime diagnosis and agent-level adaptation to LLM-based MASs for software development. Cemri *et al.* [4] selected seven open-source LLM-based MASs, and collected over 200 dialogues and execution traces, which contain the complete interaction records among all agents within the same MAS. Their analysis results show that MAS failures arise not only from limitations of the integrated

LLMs but also from the design of MASs. They finally proposed a two-tier taxonomy of failure modes in LLM-based MASs. Sarkar *et al.* [21] explored the enhancement of communication reliability and scalability in LLM-based MASs through the application of classical design patterns, with a specific focus on the Model Context Protocol (MCP). The study delved into the transition from single-agent to MASs, identifying key communication challenges such as architectural ambiguity, coordination misalignment, and task validation. Bouzenia *et al.* [2] conducted a focused empirical analysis of thought-action-result workflows produced by three contemporary LLM-based MASs. They focused on a manually annotated corpus of 120 sampled workflows unified into canonical sequences. The authors parsed heterogeneous logs, computed trajectory-level metrics, mined frequent action sub-sequences, and performed open coding of semantic relations among thoughts, actions, and results. They concluded that detecting trajectory "smells" and enforcing alignment checks or self-critique loops are promising directions to improve agent robustness.

## 2.3 Surveys on LLM-based MASs

Researchers have also conducted surveys focusing on the architecture, capabilities, and limitations of LLM-based MASs. Liu *et al.* [18] conducted a systematic literature review to investigate the architectural challenges of designing foundation model-based agents. Based on the review of 57 selected academic and industrial sources between 2023 and 2024, the authors identified 18 reusable architectural patterns that support key agent capabilities such as goal-seeking, planning, explainability, and accountability. These patterns were synthesized into a structured design pattern catalog, accompanied by a decision model to guide practitioners in selecting appropriate patterns. Their study provides holistic design knowledge for building reliable and scalable LLM-based MASs, contributing to both research and practice in MAS architectures. Liu *et al.* [17] collected 106 papers before July 1st, 2024, using keyword-based searches and snowballing from DBLP and other sources to explore what SE tasks can be solved by LLM-based agent systems and the architectural components of these systems. They categorized the works from both SE and agent perspectives, covering tasks like requirements engineering, code generation, testing, and debugging, as well as agent architectures such as planning, memory, perception, and multi-agent coordination. Their survey identified the capabilities, design patterns, and open research challenges on LLM-based agent frameworks, offering valuable insights for the design and application of LLM-based agents in SE. Chen *et al.* [5] collected 125 papers published in top artificial intelligence conferences and some unpublished yet valuable papers from arXiv in 2023 and 2024 to provide a comprehensive and up-to-date review of LLM-based MASs. Their research categorized LLM-MAS applications into three domains: solving complex tasks, simulating specific scenarios, and evaluating generative agents. For each category, the survey discussed representative systems, open-source resources, and evaluation benchmarks. They also identified key challenges for the development of LLM-based MASs, including hallucinations, limited long-context handling, communication inefficiency, and a lack of objective evaluation metrics. Yan *et al.* [27] provided a comprehensive review of LLM-based MASs by focusing on the pivotal role of communication between agents. They collected 68 papers before May 2025, including domain-specific research papers and technical papers on communication protocols, agent behaviors, and applications of LLM-based MASs in different areas. They introduced a structured framework that integrates system-level communication with system-internal communication to explore how agents interact, negotiate, and achieve collective intelligence. Yu *et al.* [28] conducted a systematic survey to investigate trustworthiness in LLM-based agents, selecting 175 representative papers from top AI, security, and NLP conferences and journals published since 2023. The authors introduced a unified conceptual framework that categorizes threats and countermeasures along two dimensions: internal and external trustworthiness. He *et al.* [10] investigated 71 papers before November 14th, 2024, indexed in DBLP, and conducted a systematic literature review complemented

by empirical case studies. They assessed and characterized the current state and potential of LLM-based MASs for SE tasks, and proposed a comprehensive research agenda aimed at strengthening individual agent competence and optimizing inter-agent collaboration.

## 2.4 Conclusive Summary

With the advances in the reasoning and decision-making capabilities of LLMs, plenty of LLM-based MASs have been designed to solve specific SE tasks. These MASs have been applied to various SE activities, such as requirements analysis, architecture design, code generation, and testing, demonstrating strong potential in simulating collaborative workflows, enabling role specialization, and enhancing automation and decision-making in software development. In addition, existing studies and surveys have investigated the design principles, agent architectures, communication mechanisms, and evaluation strategies of LLM-based MASs. These efforts provide foundational insights for constructing more robust and generalizable MAS frameworks for SE applications.

Although various studies have been conducted to investigate LLM-based MASs, no dedicated research has explored how to design an effective LLM-based MAS for a specific SE task, which QAs should be considered, which design patterns can be employed, and the underlying rationale guiding the design of LLM-based MASs for SE tasks. Our study aims to support designers in balancing QAs, selecting suitable design patterns, and explaining the design with underlying rationale when building LLM-based MASs for various SE tasks, which act as guidance for the construction of reliable, maintainable, scalable, and safe LLM-based MASs.

## 3 Study Design

In this section, we present the goal and Research Questions (RQs) (Section 3.1) formulated to investigate the design of LLM-based MASs for SE tasks. Besides, we describe the process and criteria for data collection (Section 3.2) and detail the procedures for data extraction (Section 3.3) and data analysis (Section 3.4).

### 3.1 Research Questions

The **goal** of this study is to understand how LLM-based MASs for SE tasks are designed. To this end, we formulated four Research Questions (RQs): the SE tasks addressed by LLM-based MASs (RQ1), the quality attributes prioritized by their designers (RQ2), and the design patterns (RQ3) and rationale (RQ4) guiding their construction. The four RQs are addressed by following the research process illustrated in Figure 1.
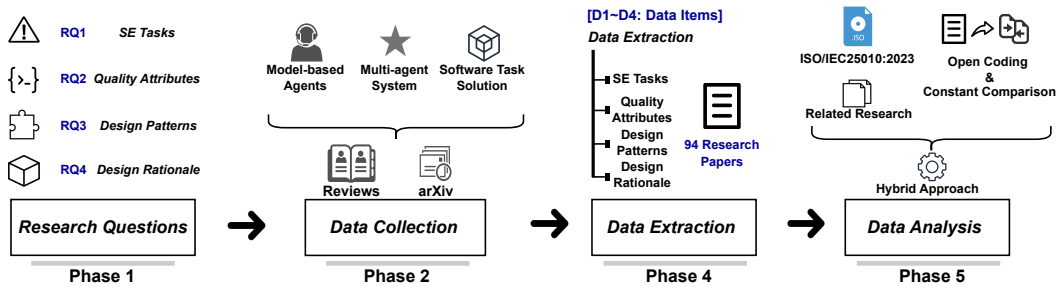


Fig. 1. Overview of the research process

0:6Cai et al.

> **RQ1: What SE tasks are addressed by LLM-based MASs?**

**Rationale**: As mentioned in Section 1, LLM-based MASs have demonstrated significant potential in addressing various SE tasks. The answer to this RQ can provide an overall view of the SE tasks addressed by LLM-based MASs, thereby informing and guiding the direction of future research in this emerging domain.

> **RQ2: What quality attributes are considered when designing LLM-based MASs for SE tasks?**

**Rationale**: Through this RQ, we aim to explore the quality attributes that designers of LLM-based MASs for SE tasks prioritize, thereby identifying the specific considerations that should be handled when designing such systems for different SE tasks.

> **RQ3: What design patterns are used for building LLM-based MASs for SE tasks?**

**Rationale**: Design patterns are reusable solutions that help balance the quality attributes of MASs. In this study, we identified the design patterns employed in building LLM-based MASs. By analyzing these patterns, we aim to uncover effective approaches that can be adopted when designing LLM-based MASs for SE tasks.

> **RQ4: What is the design rationale for constructing LLM-based MASs for SE tasks?**

**Rationale**: This RQ aims to understand the rationale behind designing LLM-based MASs for SE tasks. By examining these rationales, we seek to reveal the guiding principles that designers can follow when developing LLM-based MASs to support different SE tasks.

### 3.2 Data Collection

To obtain reliable insights into the design of LLM-based MASs for SE tasks, we relied on academic papers as our main data source. We started our data collection from the papers listed in the two recent surveys by Liu *et al.* [17] and Wang *et al.* [26], which collect and review LLM-based agent systems for SE tasks. We obtained 118 papers from the survey of Liu *et al.* [17] and 115 papers from the survey of Wang *et al.* [26]. To collect relevant papers as comprehensively as possible, we additionally performed a keyword search in the SE category on arXiv [25] using the query ("large language model" AND "agent"), retrieving papers published before 30 September 2024, when we started this study. We retained 194 papers from arXiv. After excluding duplicate papers from the three sources (i.e., [17], [26], and arXiv), we obtained a total of 236 papers. We set the following criteria to select relevant papers for our research:

(1) The paper must introduce at least one multi-agent system.
(2) The agent(s) introduced in the paper must leverage LLMs to implement their functions.
(3) The agent(s) introduced in the paper must have addressed at least one SE task.

We finally obtained 94 papers (listed in Table 7 in the Appendix: Included Studies) that served as the source for data extraction and analysis. The information from these selected papers was recorded in an MS file [3].

## 3.3 Data Extraction

To answer the four RQs presented in Section 3.1, we defined a set of data items for data extraction, as shown in Table 1. The data items D1~D4 are respectively used to extract information regarding SE tasks, QAs, design patterns, and design rationale, corresponding to RQ1~RQ4. These data items can be extracted from any section of the collected papers in our dataset, including the title, abstract, and main text. To ensure the completeness of the extracted information, we manually performed data extraction from the included papers.

Table 1. Data items extracted and their corresponding RQs

| # | Data Item | Description | RQ |
|---|---|---|---|
| D1 | SE Task | *The SE task addressed by the LLM-based MAS proposed in a paper.* | RQ1 |
| D2 | Quality Attribute | *The quality attribute(s) considered in the design of the LLM-based MAS proposed in a paper.* | RQ2 |
| D3 | Design Pattern | *The design pattern(s) employed in the design of the LLM-based MAS proposed in a paper.* | RQ3 |
| D4 | Design Rationale | *The design rationale that supports the design of the LLM-based MAS proposed in a paper.* | RQ4 |

*3.3.1 Pilot Data Extraction.* To examine whether all data items could be correctly extracted from the papers included in our dataset, we conducted a pilot data extraction. We randomly selected five papers from our dataset to perform pilot data extraction. The first author carried out the extraction independently and discussed the extraction results with the second and third authors. The results indicated that the four data items listed in Table 1 can be extracted from our dataset. Based on the results obtained from the pilot data extraction, we reached a consensus and established the following rules for the formal data extraction: (1) Only one *SE Task* can be extracted from a paper. If multiple SE tasks are addressed by the LLM-based MAS, we record the main SE task. (2) If multiple *Quality Attributes* are considered in the design of the LLM-based MAS in a paper, we record all the *Quality Attributes*. (3) If multiple *Design Patterns* are employed in the design of the LLM-based MAS in a paper, we record all the *Design Patterns*. (4) If more than one *Design Rationale* supports the design of the LLM-based MAS in a paper, we record all the *Design Rationale*.

*3.3.2 Formal Data Extraction.* After the pilot extraction, the first author independently conducted the formal data extraction based on the data items listed in Table 1. If any issues arose during the formal data extraction process, the first author discussed these issues with the second and third authors to resolve the problems. Once the first author completed the formal data extraction, the results of the extraction were reviewed by the second and third authors. Any inconsistencies were discussed among the first three authors to reach a consensus. The three authors conducted multiple rounds of reviews and revisions on the formal data extraction results to obtain the final results. The formal data extraction results were recorded in an MS Excel file and a MAXQDA file, which are publicly available in our dataset [3].

## 3.4 Data Analysis

After completing the data extraction, we conducted data analysis to answer the four RQs in Section 3.1. We categorized the *SE tasks* to answer RQ1 based on the software development lifecycle. To answer RQ2, we categorized the considered *quality attributes* based on ISO/IEC 25010:2023 [13]. To address RQ3, we used a hybrid approach to classify the employed *design patterns*. We first categorized the *design patterns* according to the category of architecture patterns for LLM-based agents provided by Liu *et al.* [18]. When no matching architecture pattern existed in that category, we subsequently employed the Open Coding and Constant Comparison methods for categorization, which are commonly used for qualitative data analysis in SE studies [24]. Additionally, we also

employed the Open Coding and Constant Comparison methods to get the category of *design rationale* for answering RQ4.

The steps we conducted for data analysis are as follows: (1) The first author read the full text of each paper in the dataset and identified the *SE task* the LLM-based MAS aiming to address, the *quality attributes* considered in the design of the LLM-based MAS, the *design patterns* employed in the design of the LLM-based MAS, and the *design rationale* supporting the design of the LLM-based MAS for SE tasks. (2) For the data items that used the Open Coding and Constant Comparison methods, the first author compared the descriptions to identify their similarities and differences, then grouped similar descriptions following a bottom-up approach, which led to the formation of higher-level categories. (3) Whenever uncertainties about the code descriptions arose, the first author discussed these uncertainties with the second and third authors to reach a consensus. Due to the iterative nature of Constant Comparison, the final results were determined after several rounds of refinement and revision. (4) The second and third authors reviewed and verified the preliminary analysis results. If any questions or discrepancies arose, the first three authors discussed them to reach a final consensus and resolve the conflicts.

## 4 Results

In this section, we report the results of the four RQs formulated in Section 3.1. For the taxonomies of *SE Task*, *Design Pattern*, *Design Rationale*, we provide a one-tier categorization. For the taxonomy of *Quality Attribute*, we provide a two-tier categorization according to ISO/IEC 25010:2023 [13].

### 4.1 Category of SE Tasks (RQ1)

Figure 2 presents the taxonomy of SE tasks extracted from our dataset [3]. It can be observed that *Code Generation* (47.87%) is the most common SE task that LLM-based MASs are designed to address. In addition, some LLM-based MASs are used to solve *Fault Localization* (9.57%), *End-to-End Software Maintenance* (8.51%), and *Program Repair* (8.51%). The remaining LLM-based MASs are identified to address *End-to-End Software Development* (7.45%), *Code Review* (6.38%), *Software Testing* (6.38%), *Requirements Engineering* (3.19%), *Code Translation* (1.06%), and *Software Deployment* (1.06%).
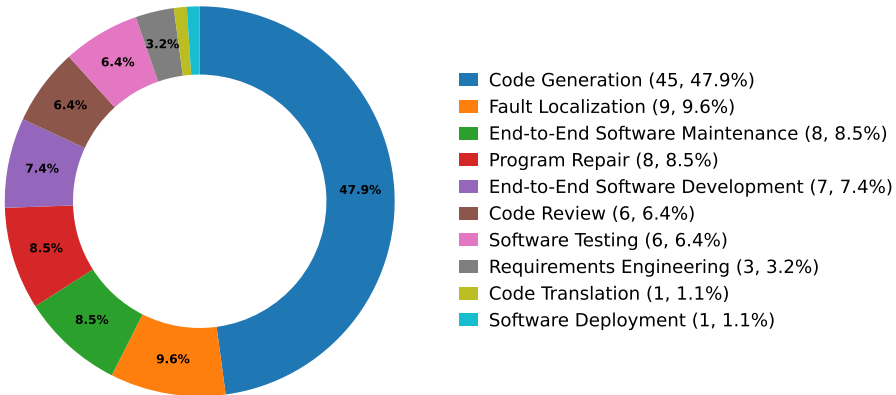


Fig. 2. Taxonomy of SE Tasks addressed by LLM-based MASs

***Code Generation*** (45, 47.9%) refers to the automated process of producing executable code, such as source code, using LLM-based MASs. Because programming is structured and rule-based, and abundant public datasets are available, many designers can use LLM-based agents to construct and

Table 2. SE Tasks addressed by LLM-based MASs

| SE Tasks | Example | Count(%) | Studies |
|---|---|---|---|
| Code Generation | *To enable LLMs to handle these real-world repo-level code generation, we present CODEAGENT, a novel LLM-based agent framework that employs external tools for effective repo-level code generation.* [S12] | 45 (47.9%) | [S1], [S7], [S9], [S11], [S12], [S15], [S20], [S23], [S24], [S27], [S29], [S33], [S34], [S36], [S37], [S38], [S40], [S41], [S42], [S45], [S46], [S51], [S52], [S53], [S54], [S58], [S61], [S62], [S63], [S64], [S67], [S68], [S71], [S73], [S74], [S76], [S78], [S79], [S82], [S83], [S84], [S86], [S87], [S92], [S94] |
| Fault Localization | *To address the limitation, this paper presents AGENTFL, a multi-agent system based on ChatGPT for automated fault localization.* [S4] | 9 (9.6%) | [S4], [S16], [S19], [S25], [S26], [S35], [S49], [S75], [S77] |
| End-to-End Software Maintenance | *In this paper, we introduce MarsCode Agent, a novel framework that leverages LLMs to automatically identify and repair bugs in software code. MarsCode Agent combines the power of LLMs with advanced code analysis techniques to accurately localize faults and generate patches.* [S31] | 8 (8.5%) | [S2], [S3], [S28], [S31], [S32], [S39], [S59], [S60] |
| Program Repair | *In this paper, we propose an automated approach for solving GitHub issues to autonomously achieve program improvement.* [S8] | 8 (8.5%) | [S8], [S14], [S17], [S50], [S69], [S70], [S89], [S90] |
| End-to-End Software Development | *MetaGPT encodes Standardized Operating Procedures (SOPs) into prompt sequences for more streamlined workflows, thus allowing agents with human-like domain expertise to verify intermediate results and reduce errors.* [S81] | 7 (7.4%) | [S5], [S13], [S21], [S57], [S65], [S81], [S93] |
| Code Review | *In this paper, we present a novel approach to improving software quality and efficiency through a Large Language Model (LLM)-based model designed to review code and identify potential issues.* [S6] | 6 (6.4%) | [S6], [S10], [S44], [S55], [S72], [S88] |
| Software Testing | *The proposed system mainly consists of three LLM-based agents responsible for action planning, state checking and parameter selecting, respectively, and two additional modules for state sensing and case rewriting.* [S43] | 6 (6.4%) | [S43], [S47], [S66], [S80], [S85], [S91] |
| Requirements Engineering | *These agents engage in product experience scenarios, through explaining their actions, observations, and challenges. Subsequent agent interviews and analysis uncover valuable user needs, including latent ones.* [S18] | 3 (3.2%) | [S18], [S30],[S48] |
| Code Translation | *In this work, we propose a novel LLM-based multi-agent system TRANSAGENT, which enhances LLM-based code translation by fixing the syntax errors and semantic errors with the synergy between four LLM-based agents...* [S56] | 1 (1.1%) | [S56] |
| Software Deployment | *In this paper, we propose GoNoGo, an LLM agent system designed to streamline automotive software deployment while meeting both functional requirements and practical industrial constraints.* [S22] | 1 (1.1%) | [S22] |

evaluate LLM-based MASs, so that *Code Generation* is the most common SE task solved by these systems. For example, Zhao *et al.* proposed an LLM-based MAS named "*VisionCoder*" that "*offers a cost-effective and efficient solution for code generation*" by "*establishing a tree-structured thought distribution and development mechanism*" [S51].

***Fault Localization*** (9, 9.6%) refers to the process of identifying the specific parts of a software system that are responsible for observing failures and errors by LLM-based MASs. By leveraging the reasoning capabilities of LLMs, designers can build MASs that assist and enhance fault localization. For instance, Wei *et al.* proposed an LLM-based MAS named "*LLM-SmartAudit*" that "*employs a collaborative system with specialized agents*" to "*detect and analyze vulnerabilities in smart contracts*" [S26].

***End-to-End Software Maintenance*** (8, 8.5%) refers to the systematic process of maintaining a software system throughout its entire workflow, from issue detection to code modification, with

the assistance of LLM-based MASs. By delegating responsibilities to specialized agents (e.g., fault localization, patch synthesis, test generation, and verification) designers of LLM-based MASs can construct automated, modular, and scalable end-to-end maintenance workflows. For example, Liu *et al.* proposed an LLM-based MAS named "*MarsCode Agent*" that "*leverages LLMs to automatically identify and repair bugs in software code*" [S31].

**Program Repair** (8, 8.5%) refers to the automated process of identifying and correcting errors and defects in software systems to restore or enhance their functionality using LLM-based MASs. By coordinating specialized agents to interpret code, which has strict syntactic and semantic constraints, program repair can be executed automatically to improve repair quality. For example, Moon *et al.* proposed "*COFFEEPOTS*", an LLM-based MAS capable of automated bug fixing through "*feedback-driven Preference-Optimized Tuning and Selection*" [S69].

**End-to-End Software Development** (7, 7.4%) refers to the entire process of building a software system using LLM-based MASs. By assigning specialized agents to distinct responsibilities (e.g., requirements elicitation, software design, implementation, and testing) designers of LLM-based MASs can construct a complete, end-to-end software development pipeline. For example, Hong *et al.* proposed an LLM-based MAS named "*MetaGPT*" that "*encodes Standardized Operating Procedures (SOPs) into prompt sequences for more streamlined workflows*" to allow the system "*with human-like domain expertise*" in software development [S81].

**Code Review** (6, 6.4%) refers to the systematic process of examining source code performed by LLM-based MASs to identify defects, ensure compliance with coding standards, and improve code quality. By employing specialized agents for automated style and correctness checking, semantic analysis, and reviewer suggestion generation, LLM-based MASs can deliver consistent and context-aware feedback that accelerates the review process. For instance, Tang *et al.* proposed an LLM-based MAS named "*CodeAgent*" which is aimed at "*for code review automation*" [S44].

**Software Testing** (6, 6.4%) refers to the automated process of evaluating a software system and its components to verify whether they meet specified requirements and to identify defects and deviations from expected behavior by using LLM-based MASs. By coordinating LLM-based agents for test-case generation, prioritization, oracle creation, and automated execution, MASs can provide adaptive and context-aware test suites that increase defect detection and testing efficiency. For instance, Yoon *et al.* proposed an LLM-based MAS named "*DROIDAGENT*", which is "*an autonomous GUI testing agent for Android, for semantic, intent-driven automation of GUI testing*" [S47].

**Requirements Engineering** (3, 3.2%) refers to the comprehensive process of analyzing, specifying, validating, and managing the requirements of a software system by leveraging LLM-based MASs. By assigning agents to tasks such as elicitation, conflict detection, prioritization, and validation, designers can build an automated, collaborative, and traceable workflow for requirements engineering. For example, Jin *et al.* proposed an LLM-based MAS named "*MARE*" that "*leverages collaboration among LLMs throughout the entire RE process*" to obtain "*high-quality requirements specifications*" [S30].

**Code Translation** (1, 1.1%) refers to the systematic process of transforming source code written in one programming language into equivalent code in another programming language, while preserving the original functionality and behavioral integrity with the assistance of LLM-based MASs. By leveraging the language understanding and reasoning capabilities of LLMs, code translation can be achieved through mapping and transforming the syntactic and semantic rules across

different programming languages. For example, Yuan *et al.* proposed an LLM-based MAS named "*TRANSAGENT*" that "*leverages parallel data to train models for automated code translation*" [S56].

***Software Deployment*** (1, 1.1%) refers to the systematic process of delivering, installing, and configuring a software system and its updates in a target environment to make it operational for end users by using LLM-based MASs. By coordinating agents responsible for environment configuration, dependency management, and monitoring, MASs can support automated and reliable deployment processes that minimize risks and reduce operational overhead. For example, Khoee *et al.* proposed an LLM-based MAS named "*GoNoGo*" which is "*designed to streamline automotive software deployment while meeting both functional requirements and practical industrial constraints*" [S22].

## 4.2   Category of Quality Attributes (RQ2)

As mentioned in Section 3.1, the LLM-based MASs for SE tasks involve various quality attributes that designers consider during their construction and implementation. In this section, we show the results of RQ2 according to ISO/IEC 25010:2023 [13]. Figure 3 and Table 3 present the results of data extraction and analysis for RQ2. Results show that *Functional Suitability* (94.7%) is the QA that is mostly considered in the design of LLM-based MASs to address SE tasks. Meanwhile, *Performance Efficiency* (51.1%) and *Maintainability* (50.0%) are also considered by a number of LLM-based MASs. The remaining QAs are *Reliability* (36.2%), *Flexibility* (25.5%), *Compatibility* (10.6%), *Security* (10.6%), and *Interaction Capability* (9.6%).
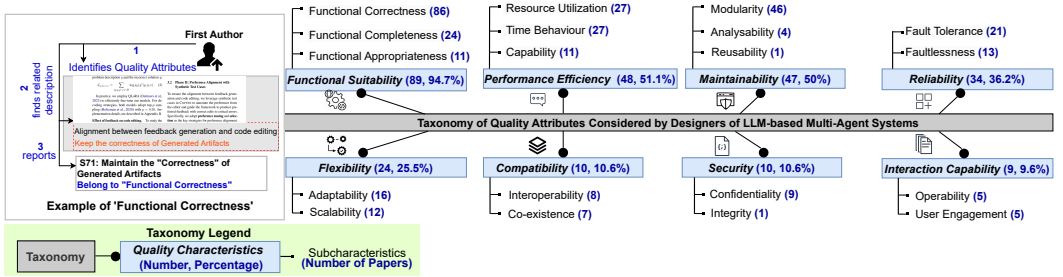


Fig. 3. Taxonomy of Quality Attributes considered in the design of LLM-based MASs for SE Tasks

***Functional Suitability*** (89, 94.7%) is the ability of an LLM-based MAS to provide functions that meet both explicit and implicit requirements of its intended users when operated under designated conditions. In the data extraction results, there are three attributes considered by the designers of LLM-based MASs in the category of *Functional Suitability*:

- *Functional Correctness* (86): The capacity of an LLM-based MAS to deliver a range of functions that encompasses all specified tasks and objectives of intended users.
- *Functional Completeness* (24): The ability of an LLM-based MAS to yield precise results when utilized by intended users.
- *Functional Appropriateness* (11): The capacity of an LLM-based MAS to deliver functions that support the achievement of specified tasks and objectives. It indicates that ensuring a system satisfies specified functional requirements and produces correct outputs remains the primary concern for most designers of LLM-based MASs for SE tasks. *Functional Suitability*, especially *Functional Correctness*, is critically important for the QAs of LLM-based MASs. For example, Zhang *et al.* proposed a method called "*Deep Retrieval-Augmented Generation*" in

Table 3. Quality Attributes considered in the design of LLM-based MASs for SE Tasks

| Quality Attribute | Sub_Quality Attribute | Studies |
|---|---|---|
| Functional Suitability | Functional Correctness | [S1], [S2], [S3], [S4], [S5], [S7], [S8], [S9], [S11], [S12], [S13], [S14], [S15], [S16], [S17], [S19], [S20], [S22], [S23], [S24], [S25], [S26], [S27], [S28], [S29], [S30], [S31], [S32], [S33], [S34], [S36], [S38], [S39], [S40], [S41], [S42], [S43], [S44], [S46], [S47], [S48], [S49], [S50], [S51], [S52], [S53], [S54], [S55], [S56], [S57], [S58], [S59], [S61], [S62], [S63], [S64], [S65], [S66], [S67], [S68], [S69], [S70], [S71], [S72], [S73], [S74], [S75], [S76], [S77], [S78], [S79], [S80], [S81], [S82], [S83], [S84], [S85], [S86], [S87], [S88], [S89], [S90], [S91], [S92], [S93], [S94] |
| | Functional Completeness | [S6], [S12], [S18], [S21], [S24], [S29], [S30], [S34], [S44], [S47], [S54], [S55], [S62], [S64], [S73], [S76], [S79], [S80], [S83], [S84], [S85], [S91], [S93], [S94] |
| | Functional Appropriateness | [S22], [S34], [S36], [S41], [S49], [S55], [S67], [S69], [S73], [S93], [S94] |
| Performance Efficiency | Resource Utilization | [S2], [S4], [S6], [S8], [S18], [S19], [S20], [S25], [S28], [S35], [S36], [S37], [S39], [S40], [S45], [S46], [S47], [S48], [S49], [S58], [S63], [S64], [S68], [S81], [S87], [S88], [S93] |
| | Time Behaviour | [S2], [S4], [S6], [S8], [S10], [S13], [S21], [S22], [S27], [S39], [S45], [S46], [S47], [S49], [S55], [S56], [S58], [S60], [S65], [S73], [S80], [S81], [S82], [S87], [S90], [S91], [S94] |
| | Capability | [S18], [S26], [S38], [S45], [S48], [S57], [S59], [S85], [S88], [S89], [S94] |
| Maintainability | Modularity | [S2], [S4], [S5], [S6], [S7], [S11], [S14], [S15], [S16], [S18], [S21], [S26], [S29], [S30], [S31], [S32], [S36], [S38], [S43], [S44], [S47], [S49], [S51], [S55], [S56], [S59], [S60], [S61], [S62], [S63], [S64], [S66], [S67], [S68], [S69], [S70], [S71], [S72], [S74], [S76], [S78], [S81], [S82], [S84], [S88], [S94] |
| | Analysability | [S6], [S67], [S92], [S94] |
| | Reusability | [S74] |
| Reliability | Fault Tolerance | [S5], [S7], [S8], [S10], [S30], [S39], [S40], [S42], [S47], [S55], [S57], [S61], [S63], [S65], [S70], [S73], [S81], [S83], [S84], [S85], [S87], [S88], [S92] |
| | Faultlessness | [S1], [S13], [S24], [S27], [S48], [S51], [S52], [S53], [S62], [S64], [S77], [S86], [S94] |
| Flexibility | Adaptability | [S1], [S2], [S9], [S12], [S31], [S33], [S36], [S52], [S53], [S57], [S66], [S75], [S81], [S82], [S88], [S91] |
| | Scalability | [S2], [S13], [S15], [S17], [S33], [S37], [S38], [S52], [S69], [S82], [S89], [S94] |
| Compatibility | Interoperability | [S12], [S23], [S25], [S46], [S55], [S76], [S79], [S84] |
| | Co-Existence | [S2], [S12], [S23], [S25], [S76], [S79], [S80] |
| Security | Confidentiality | [S3], [S9], [S21], [S25], [S44], [S65], [S77], [S82], [S85] |
| | Integrity | [S64] |
| Interaction Capability | User Engagement | [S20], [S21], [S40], [S74], [S84] |
| | Operability | [S40], [S46], [S49], [S66], [S76] |

their LLM-based MAS to "*handle the complex inheritance relationships in exception types*"; thus, they can ensure the correctness of exception handling [S55].

**Performance Efficiency** (48, 51.1%) is the ability of an LLM-based MAS to execute its functions within the designated time and performance parameters while maintaining efficient use of resources (e.g., CPU, memory, storage) under specified conditions. In the data extraction results, there are three attributes considered by the designers of LLM-based MASs in the category of *Performance Efficiency*:

- *Resource Utilization* (27): The ability of an LLM-based MAS to perform designated functions under specified conditions while consuming no more than the allocated amount of resources.
- *Time Behaviour* (27): The capacity of an LLM-based MAS to execute designated functions under specified conditions, ensuring that response time and throughput rates satisfy established requirements.
- *Capability* (11): The capacity of an LLM-based MAS to fulfill requirements for maximum limits of system parameters (e.g., concurrent user support). It indicates that many designers of LLM-based MASs prioritize optimizing computational resource use and ensuring timely system responses, while the systems can give responses of high functional quality. For example, Zhang *et al.* let the MAS "*whether the information is parameterized*" to "introduces minimal latency at inference time " as the extraneous information has been eliminated [S45].

***Maintainability*** (47, 50.0%) is the capacity of an LLM-based MAS to be effectively and efficiently modified by designated maintainers, such as corrections, enhancements, and adaptations to environmental changes. In the data extraction results, there are three attributes considered by the designers of LLM-based MASs in the category of *Maintainability*:

- *Modularity* (46): The ability of an LLM-based MAS to confine modifications to a single component, minimizing impact on other components, with the system structured into discrete modules and components exhibiting high cohesion and minimal coupling to other modules and components.
- *Analysability* (4): The capacity of an LLM-based MAS to be effectively and efficiently evaluated for the impact of proposed changes to one or more components, to diagnose deficiencies, failure causes, and to identify components requiring modification, incorporating mechanisms for self-analysis of faults and generation of reports prior to failures and other events.
- *Reusability* (1): The capacity of an LLM-based MAS to be utilized as a component in multiple systems and in the development of other assets. It indicates that various designers of LLM-based MASs are willing to build modular architectures to simplify maintenance, enable isolated updates, debugging, and evolution in these systems. For example, Zhao *et al.* introduced a special agent called "*Team Leader*" to "*divide the project into distinct branch modules and define the functional scope and objectives of each module*". In this way, they can improve the modularity of the MAS [S51].

***Reliability*** (34, 36.2%) is the ability of an LLM-based MAS to execute designated functions under specified conditions for a defined duration without interruptions and failures. In the data extraction results, there are two attributes considered by the designers of LLM-based MASs in the category of *Reliability*:

- *Fault Tolerance* (21): The capability for an LLM-based MAS to function as intended even in the presence of hardware and software faults.
- *Faultlessness* (13): The capability of an LLM-based MAS to execute defined functions without faults during normal operation. In addition, the conception of *Faultlessness* may also be extended to other QAs to reflect the extent to which they fulfill required needs under typical operating conditions. It indicates that plenty of designers of LLM-based MASs place strong importance on ensuring system robustness to reduce potential failures and unexpected errors. For example, Zhang *et al.* used a "*retry-loop of generating patches*" in their LLM-based MAS, "*if a generated patch does not follow the specified patch format or cannot be applied syntactically to the original program, the agent is prompted to retry*" [S8].

***Flexibility*** (24, 25.5%) is the capability of an LLM-based MAS to adjust in response to changes in requirements, usage contexts, and the system environment. Flexibility with respect to usage context involves two distinct aspects: the technical (e.g., software environment) and the non-technical (e.g., physical environment). There are two attributes considered by the designers of LLM-based MASs in the category of *Flexibility*:

- *Adaptability* (16): The capability of an LLM-based MAS to be effectively and efficiently modified and transferred for operation in different hardware, software, and operational environments.
- *Scalability* (12): The capability of an LLM-based MAS to manage increasing or decreasing workloads, and to adjust its capacity in response to fluctuations in requirements. It indicates that some designers of LLM-based MASs focus on enabling systems to adjust to evolving requirements, dynamic runtime environments, and diverse application contexts. For example,

Hu *et al.* design an LLM-based MAS called "*EvoMAC*", which has "*ability to iteratively adapt both agents and their connections during test time for each task*" [S53].

**Compatibility** (10, 10.6%) is the capability of an LLM-based MAS to exchange information with other systems and to carry out required functions while operating within a shared environment and utilizing common resources. In the data extraction results, there are two attributes considered by the designers of LLM-based MASs in the category of *Compatibility*:

- *Interoperability* (8): The capability of an LLM-based MAS to exchange information with other systems and to make mutual use of the exchanged information.
- *Co-Existence* (7): The capability of an LLM-based MAS to perform required functions efficiently within a shared environment and using common resources, without causing negative effects on the operation of other systems. It indicates that some designers of LLM-based MASs place importance on enabling seamless interaction and integration with other systems, platforms, and tools. For example, Zhang *et al.* proposed an LLM-based MAS called "*CODEAGENT*" which "*integrates five programming tools, enabling interaction with software artifacts for information retrieval, code implementation, and code testing*" [S12].

**Security** (10, 10.6%) is the capability of an LLM-based MAS to safeguard information and data by ensuring that individuals or other systems access data according to their designated types and authorization levels, and by resisting attacks from malicious entities, which includes protection of both stored data and data in transmission. In the data extraction results, there are two attributes considered by the designers of LLM-based MASs in the category of *Security*:

- *Confidentiality* (9): The capability of an LLM-based MAS to ensure that data can be accessed exclusively by entities with appropriate authorization.
- *Integrity* (1): The capacity of an LLM-based MAS to safeguard the integrity of its system and data against unauthorized alteration and deletion, whether due to malicious intent or computational error. It indicates that some designers of LLM-based MASs specially focus on protecting sensitive data and ensuring that information is accessible only to authorized entities. For example, Talebirad *et al.* used "*a stateless oracle agent, which can monitor each sensitive task and decide if it is indeed malicious or not*" to promise the confidentiality of their LLM-based MAS [S82].

**Interaction Capability** (9, 9.6%) is the ability of an LLM-based MAS that enables specified users to exchange information through the user interface for achieving intended tasks. In the data extraction results, there are two attributes considered by the designers of LLM-based MASs in the category of *Interaction Capability*:

- *Operability* (5): The capability of an LLM-based MAS to incorporate functions and attributes that facilitate ease of operation and control, which is closely associated with controllability, robustness against user errors, and alignment with user expectations.
- *User Engagement* (5): The ability of an LLM-based MAS to present functions and information in a manner that is appealing and motivating, thereby encouraging sustained user interaction, which encompasses system properties that enhance user pleasure and satisfaction, including harmonious color schemes, intuitive interface design, and friendly voice-based guidance. It indicates that some designers of LLM-based MASs pay special attention to ensuring ease of operation and effective human–system interaction to improve usability and user experience. For example, Josifosk *et al.* proposed an LLM-based MAS called "*Flows*", which "*supports research in the design of interactions involving humans as computational building blocks in a*

*way that maximizes the utility of the overall computation with minimal human effort*". This method promises that the whole system is friendly to human beings [S74].

## 4.3 Category of Design Patterns (RQ3)

As mentioned in Section 3.1, numerous design patterns are employed to guide the construction and implementation of LLM-based MASs, and some systems have multiple design patterns concurrently. In this section, we report the results of RQ3 under the guidance of the work of Liu *et al.* [18]. Table 4 presents the taxonomy of the design patterns extracted from the dataset. Results show that *Role-Based Cooperation* (46.8%) is the most frequently used design pattern by the designers of LLM-based MASs. *Self-Reflection* (36.2%) is also widely used by the designers of LLM-based MASs for SE tasks. The remaining design patterns are *Tool-Agent Registry* (14.9%), *Cross-Reflection* (12.8%), *Retrieval-Augmented Generation (RAG)* (10.6%), *Human-Reflection* (6.4%), *Agent Adapter* (5.3%), *Single-Path Plan Generator* (5.3%), *Prompt/Response Optimiser* (4.3%), *Debate-Based Cooperation* (4.3%), *Layered-Based Cooperation* (4.3%), *Agent Evaluator* (3.2%), *Multi-Path Plan Generator* (3.2%), *Incremental Model Querying* (3.2%), *Hierarchical Coordination* (3.2%), and *Voting-Based Cooperation* (3.2%).

**Role-Based Cooperation** (44, 47.4%) refers to a collaborative pattern in which each LLM-based agent assumes a distinct role that defines its functional responsibilities and interaction protocols within the system. These roles are intended to optimize the overall performance of the agents by utilizing their specialized capabilities and ensuring efficient task allocation and coordination. Employing *Role-Based Cooperation* in the design of MASs facilitates clear responsibility allocation, enhances scalability and adaptability, and supports flexible coordination strategies. For example, Dong *et al.* proposed an LLM-based MAS, and divided the agent group into three roles: "*Analyst, Coder, Tester*" [S86].

**Self-Reflection** (34, 35.8%) refers to a pattern in which each agent uses the capabilities of LLMs to perform introspective analysis of its own actions, decisions, and performance. Through this process, the agent can identify areas for improvement, adapt to dynamic environments, and optimize its contributions to collaborative tasks, thereby enhancing the overall effectiveness of agent cooperation. It shows that designers of LLM-based MASs enable agents to introspect and iteratively refine their outputs, reducing errors and adapting to changing task requirements. For instance, Hong *et al.* proposed an LLM-based MAS called "*MetaGPT*"; its code agent can "*improve code using its own historical execution and debugging memory*" with the help of *Self-Reflection* [S81].

**Tool-Agent Registry** (14, 14.7%) refers to a centralized repository that maintains a comprehensive and unified resource which is readily accessible for assigning appropriate tools to agents according to task requirements to enable effective and reliable task execution. The registry serves as an integral mechanism for agent collaboration by enabling dynamic discovery, selection, and utilization of resources. Through this functionality, it improves output quality, reduces model workload and latency, thereby enhancing the flexibility of the whole system. For example, Fan *et al.* proposed an LLM-based MAS called "*ICAA*", and set up a special "*agent integrating a toolbox that includes context-aware splitting, code retrieval, documentation retrieval, Web search, and static code analysis tools*". This agent incorporates a "*thinking–decision–action loop*" to select and schedule tools for different "*sub-agents*" [S88].

**Cross-Reflection** (12, 12.6%) refers to a collaborative pattern in which each agent uses LLMs within the agent system to observe, analyze, and interpret the behaviors, decisions, and outputs of other

Table 4. Design Patterns Used by Designers of LLM-based MASs for SE Tasks

| Design Pattern | Example | Count(%) | Studies |
|---|---|---|---|
| Role-Based Cooperation | *Specifically, INTERVENOR employs two LLM-based agents to play different roles in code repair. [S78]* | 44 (46.8%) | [S2], [S4], [S5], [S7], [S11], [S13], [S14], [S16], [S18], [S19], [S21], [S26], [S28], [S29], [S30], [S31], [S32], [S35], [S36], [S37], [S38], [S42], [S44], [S46], [S47], [S49], [S51], [S55], [S56], [S57], [S58], [S59], [S60], [S61], [S63], [S64], [S67], [S69], [S71], [S78], [S81], [S86], [S93], [S94] |
| Self-Reflection | *..., in order to reduce hallucinations and inefficient planning, we apply a self-reflection mechanism. [S43]* | 34 (36.2%) | [S1], [S2], [S3], [S11], [S19], [S22], [S24], [S27], [S31], [S36], [S37], [S38], [S39], [S42], [S43], [S45], [S47], [S49], [S53], [S57], [S58], [S60], [S62], [S63], [S72], [S73], [S74], [S75], [S79], [S81], [S87], [S90], [S91], [S92] |
| Tool-Agent Registry | *Across five tasks with Python and WikiSearch API as tools, Lemur-70B-Chat outperforms both Llama-2-70B-Chat and CodeLlama-34B-INST by large margins. [S79]* | 14 (14.9%) | [S2], [S8], [S9], [S12], [S25], [S27], [S75], [S76], [S79], [S80], [S82], [S85], [S88], [S94] |
| Cross-Reflection | *Meanwhile, this Stderr information is provided to the questioner, who will generate a natural language description based on the Stderr. This natural language description is also appended to the dialogue messages. Next, both the natural language description and the Stderr are provided as new questions to the programmer, who will continue to modify the code. [S7]* | 12 (12.8%) | [S7], [S13], [S26], [S41], [S50], [S61], [S65], [S67], [S69], [S71], [S74], [S78] |
| Retrieval-Augmented Generation (RAG) | *Furthermore, we draw inspiration from the RAG (Retrieval-Augmented Generation) approach, where we match similar content from a memory pool as additional information. [S36]* | 10 (10.6%) | [S8], [S23], [S36], [S45], [S51], [S54], [S55], [S64], [S89], [S94] |
| Agent Adapter | *We give the agent access to tools, including access to: ... [S25]* | 6 (6.4%) | [S6], [S9], [S22], [S25], [S27], [S94] |
| Human-Reflection | *AISD is designed to keep the user in the loop, i.e., by repeatedly taking user feedback on use cases, high-level system designs, and prototype implementations through system testing. [S20]* | 5 (5.3%) | [S20], [S27], [S40], [S84], [S94] |
| Single-Path Plan Generator | *AXNav consists of three main components that are used to prepare for, execute, and export test results: ... [S66]* | 5 (5.3%) | [S10], [S15], [S48], [S66], [S68] |
| Prompt/Response Optimiser | *The Requirements Engineering agent is an integral part of our engine, leveraging AI capabilities to automate the generation and prioritization of software requirements. [S21]* | 4 (4.3%) | [S4], [S9], [S21], [S33] |
| Debate-Based Cooperation | *To address this, we implemented a Multi-Agent Debate (MAD) mechanism to establish a loop between generator and validator. [S3]* | 4 (4.3%) | [S3], [S35], [S73], [S83] |
| Layered-Based Cooperation | *..., we propose a layered approach for implementing capabilities in LLM-based applications by mapping them to the layers and components with corresponding attributes. [S45]* | 4 (4.3%) | [S15], [S23], [S45], [S77] |
| Agent Evaluator | *To verify that generated fix suggestions can plausibly fix each issue, FixAlly evaluates each code modification in its Suggestion Assessment module. [S60]* | 3 (3.2%) | [S60], [S61], [S70] |
| Multi-Path Plan Generator | *It divides the Vulnerability Identification phase into 40 targeted scenarios, each focused on a specific vulnerability type (as listed in our repository). [S26]* | 3 (3.2%) | [S14], [S23], [S26] |
| Incremental Model Querying | *These patches might be from running a single SWE agent multiple times or running multiple SWE agents. [S17]* | 3 (3.2%) | [S4], [S15], [S17] |
| Hierarchical Coordination | *Our method employs a novel hierarchical coordination paradigm, inspired by a cognitive debugging model, to efficiently manage cognitive steps with minimal communication and dynamically adjust to bug complexity through its three-level architecture. [S2]* | 3 (3.2%) | [S2], [S34], [S76] |
| Voting-Based Cooperation | *Finally, the model merges and votes on all candidate solutions, selecting the highest-voted one as the final repair solution. [S31]* | 3 (3.2%) | [S17], [S31], [S33] |

agents in the MAS. This process allows agents to gain insights into the strategies and performance of their peers, enabling adaptive learning, strategic alignment, and improved coordination. It also supports error detection and correction, reduces individual biases, and fosters collective calibration and consensus. For example, Li *et al.* proposed an LLM-based MAS called "*CAMEL*", which used a method called "*Critic-In-Loop*" to implement *Cross-Reflection* [S67].

***Retrieval-Augmented Generation (RAG)*** (10, 10.5%) refers to a collaborative pattern that integrates retrieval-based methods with generative capabilities to enhance the problem-solving and decision-making processes of agents. This approach allows agents to retrieve relevant information from external knowledge sources and databases and integrate it seamlessly into their generated outputs and actions. RAG enables agents to access up-to-date and contextually appropriate data, enhancing the accuracy, relevance, and effectiveness of their contributions to collaborative tasks. For instance, Wu *et al.* proposed an LLM-based MAS called "*AutoGen*" that uses "*RAG*" to enhance the quality of code generation and question answering [S64].

***Agent Adapter*** (6, 6.3%) refers to a pattern that enables smooth interaction and collaboration among different agents by standardizing communication protocols, translating data formats, and managing interfaces. This pattern allows agents to exchange information and coordinate actions effectively, thereby improving the overall performance and adaptability of the multi-agent system. For example, Rasheed *et al.* proposed an LLM-based MAS, using GitHub APIs as the *Agent Adapter* to "*systematically download and process the code*" [S6].

***Human-Reflection*** (5, 5.3%) refers to a collaborative pattern that human users and experts observe, evaluate, and provide reflective feedback on the actions, decisions, and outputs of the agents. This feedback is then utilized by the agents to adjust their strategies, improve their performance, and enhance their collaborative capability, thereby ensuring better alignment with human goals and ethical considerations. For instance, Thoppilan *et al.* proposed an LLM-based MAS called "*LaMDA*", which uses the feedback of "*other crowdworkers*" in the reflection of agents [S94].

***Single-Path Plan Generator*** (5, 5.3%) refers to a collaborative pattern tasked with designing a linear and sequential plan to accomplish a specific goal and task. Initially, a Planner agent (or a team of cooperating Planner agents) constructs a cohesive, linear, and executable plan to achieve a specified objective. The resulting plan consists of a single, non-branching pathway of actions and decisions, determined by analyzing the input data, the operational capabilities of the involved agents, and the environmental constraints. The plan is decomposed into constituent subtasks and allocates them to other agents for cooperative execution. This pattern provides a clear global plan that reduces coordination ambiguity and conflict among agents, and enables efficient task decomposition and deterministic assignment of subtasks which lowers negotiation overhead and improves throughput. For example, Taeb *et al.* proposed an LLM-based MAS named "*AXNav*", which can "*formulate a step-by-step plan that accomplishes the goal*" provided by the users of the MAS [S66].

***Prompt/Response Optimiser*** (4, 4.3%) refers to a specialized pattern that dynamically refines and enhances the prompts provided to LLMs and optimizes the generated responses to improve their relevance, accuracy, and alignment with the objectives of the agent system. This pattern employs iterative feedback loops, context analysis, and performance metrics to fine-tune input prompts and evaluate output quality, ensuring effective communication and coordination among agents.

For instance, Tufano *et al.* proposed an LLM-based MAS named "*AutoDev*", which can refine the prompts provided by the users "*in a predefined format*" [S9].

**Debate-Based Cooperation** (4, 4.3%) refers to a collaborative pattern by which multiple agents, using the reasoning and argumentative capabilities of LLM, engage in structured debates to evaluate various perspectives, challenge assumptions and converge on optimal solutions for a given task. Each agent adopts a distinct viewpoint, presenting arguments supported by evidence and logic, while critically assessing the arguments of others. This pattern enhances MAS design by structuring adversarial argumentation among agents to surface diverse perspectives, reveal and correct errors and weak reasoning. For example, Zhang *et al.* proposed an LLM-based MAS named "*ACFIX*", which introduced "*Multi-Agent Debate (MAD) mechanism*" between the different agent roles [S3].

**Layered-Based Cooperation** (4, 4.3%) refers to a structured collaborative pattern wherein agents are organized into hierarchical layers, each layer assigned specific agent roles, tasks, and levels of abstraction based on their capabilities. Agents within each layer leverage the generative and reasoning capacities of LLMs to process inputs, execute tasks, and produce outputs that are passed to subsequent layers for further refinement and integration. This pattern separates concerns across hierarchical abstraction levels, enabling role specialization, modularity, and scalable coordination. For instance, Zan *et al.* proposed an LLM-based MAS named "*CODES*", which is formulated in a "*Multi-Layer Sketch*" to accomplish the repository-level code generation [S15].

**Agent Evaluator** (3, 3.2%) refers to a specialized pattern responsible for assessing the performance, decisions, and outputs of other agents within the MAS. *Agent Evaluator* is an independent, specialized role or subsystem that systematically conducts metric-based evaluation and governance of the performance, decisions, and outputs of other agents, thereby supporting targeted remediation, quality assurance, and continuous improvement of the MAS. For example, Huang *et al.* proposed an LLM-based MAS named "*AgentCoder*", which implements a "*test agent*" to evaluate the result of code generation. "*If all test cases are passed, the agent returns the code to the human developer*" [S61].

**Multi-Path Plan Generator** (3, 3.2%) refers to a specialized pattern that generates multiple alternative paths and action options of agents at each intermediate step, then assembles them into the final task execution plan. Each plan represents a unique sequence of actions and decisions that account for various scenarios, agent capabilities, and environmental variables. This pattern increases the robustness and adaptability of an LLM-based MAS by generating multiple plans that support the selection of optimal or complementary strategies across varying agent capabilities and environmental conditions. For example, Ma *et al.* proposed an LLM-based MAS, which can generate multiple plans to complete the tasks given by users, and the MAS "*iteratively narrows down the search space and guides the agents to focus on the most relevant area by simulating multiple workflows and evaluating their reward score*" [S23].

**Incremental Model Querying** (3, 3.2%) refers to a collaborative pattern that agents iteratively query LLMs by breaking down complex tasks and queries into smaller, sequential sub-queries. Each query builds on the results of previous ones, enabling agents to refine their understanding, integrate intermediate findings, and gradually develop a comprehensive response and solution. This pattern breaks down complex tasks into sequential sub-queries that iteratively enhance understanding, incorporate intermediate results, and converge toward a high-quality solution. For example, Qin *et al.* proposed an LLM-based MAS called "*AGENTFL*", which adopts "*a Multi-Round Dialogue strategy in task*" to improve the efficiency of fault localization [S4].

***Hierarchical Coordination*** (3, 3.2%) refers to a dynamic and scalable collaborative pattern that employs a tiered approach to resolve given SE tasks. Initially, a minimal set of agents that leverage the capabilities of LLMs collaborate through simple coordination strategies to address a given task. If the task remains unsolved, the multi-agent system gradually expands the collaboration framework by introducing additional agents and tools, thereby enhancing its complexity and capacity. This iterative process continues until the problem is successfully solved or the system reaches its predefined complexity limit. For example, Lee *et al.* proposed an LLM-based MAS called "*FixAgent*" under the guidance of *Hierarchical Coordination* to improve the efficiency of problem solving [S2].

***Voting-Based Cooperation*** (3, 3.2%) refers to a collaborative pattern in which multiple agents, using the reasoning and evaluation capabilities of LLMs, make collective decisions by voting on proposed actions, strategies, and solutions. Each agent contributes its perspective based on its role, expertise, and task analysis, and the final decision is determined through a predefined voting protocol, such as majority rule or weighted voting. This pattern enables democratic participation, reduces individual bias, and promotes consensus-based outcomes, thereby enhancing the robustness and fairness of collaborative task execution. For instance, Li *et al.* proposed an LLM-based MAS, which uses the voting between agents to "*determine the final answer*" to the SE problems [S33].

## 4.4   Category of Design Rationale (RQ4)

As mentioned in Section 3.1, all these LLM-based MASs for SE tasks are developed following the design rationale that guides the construction of LLM-based MASs to support SE tasks. In this section, we report the results of RQ4. Table 5 presents the taxonomy of the design rationale extracted from the dataset. Results show that *Improving the Quality of Generated Code* (44.7%) is the most used design rationale. At the same time, *Simulating Human Processes of Solving SE Tasks* (29.8%) is also used by various designers of LLM-based MASs. The remaining design rationale: *Optimizing Software Resource Management* (28.7%), *Improving the Efficiency of Generating Software Artifacts* (24.5%), *Reducing the Difficulty of Task Resolution* (20.2%), *Improving the Adaptability of Agent System* (19.1%), *Enhancing the Diversity of Generated Software Artifacts* (12.8%), and *Ensuring Software Security* (4.3%).

***Improving the Quality of Generated Code*** (42, 44.7%) is the most frequently used design rationale. The quality of code generated by LLMs directly impacts the quality of software artifacts produced by LLM-based MASs. Consequently, lots of designers place special emphasis on improving code quality, and they employ a variety of methods to achieve this goal. For example, Olausson *et al.* proposed an LLM-based MAS, which uses "*the error message feedback generated from execution environment*" to improve the quality of generated code [S92].

***Simulating Human Processes of Solving SE Tasks*** (28, 29.8%) is also adopted by plenty of designers. Various LLM-based MASs draw design inspiration from human processes of solving SE tasks and are designed and implemented based on this foundation. For example, Zhang *et al.* proposed an LLM-based MAS named "*Self-Edit*", which "*is inspired by the problem-solving process of human programmers*" to improve the accuracy of code generation [S87].

***Optimizing Software Resource Management*** (27, 28.7%) is also a widely used rationale. Many designers of LLM-based MASs focus on the efficient allocation, utilization, and coordination of computational and memory assets among autonomous agents. For example, Qin *et al.* proposed an LLM-based MAS named "*AGENTFL*" decomposing the whole fault localization process into

Table 5.  Design Rationale of LLM-based MASs for SE Tasks

| Design Rationale | Example | Count(%) | Studies |
|---|---|---|---|
| Improving the Quality of Generated Code | *We instructed each agent to iterate with the others at least three times to refine and update the code based on feedback from the previous agents, ensuring a thorough and collaborative development process. [S13]* | 42 (44.7%) | [S1], [S5], [S6], [S7], [S11], [S12], [S13], [S24], [S27], [S29], [S31], [S33], [S36], [S40], [S41], [S42], [S45], [S46], [S47], [S49], [S50], [S51], [S52], [S54], [S57], [S58], [S59], [S60], [S61], [S62], [S64], [S65], [S66], [S70], [S73], [S79], [S81], [S83], [S90], [S91], [S92], [S93] |
| Simulating Human Processes of Solving SE Tasks | *Specifically, iAudit is inspired by the observation that expert human auditors first perceive what could be wrong and then perform a detailed analysis of the code to identify the cause. [S16]* | 28 (29.8%) | [S2], [S5], [S7], [S10], [S11], [S12], [S14], [S16], [S19], [S26], [S28], [S29], [S30], [S35], [S42], [S44], [S53], [S55], [S58], [S59], [S64], [S67], [S71], [S72], [S78], [S81], [S86], [S87] |
| Optimizing Software Resource Management | *We posit that our approach will considerably augment the field of software quality assurance, rendering it more efficient, precise, and cost-effective. [S88]* | 27 (28.7%) | [S2], [S4], [S6], [S8], [S14], [S18], [S19], [S20], [S25], [S28], [S32], [S34], [S35], [S36], [S37], [S40], [S45], [S46], [S47], [S48], [S49], [S58], [S63], [S68], [S81], [S87], [S88] |
| Improving the Efficiency of Generating Artifacts | *Moreover, our analysis of agent interactions within AGENTVERSE reveals the emergence of specific collaborative behaviors, contributing to heightened group efficiency. [S62]* | 23 (24.5%) | [S2], [S3], [S4], [S9], [S13], [S17], [S19], [S21], [S27] , [S39], [S40], [S45], [S49], [S61], [S62], [S63], [S64], [S73], [S75], [S87], [S88], [S90], [S91] |
| Reducing the Difficulty of Task Resolution | *Specifically, by decomposing the localization process into several steps and achieving FL through the collaboration of multiple LLM-driven agents (i.e., intelligent entities that can perceive the environment, make decisions, and perform actions), the advantages we expect are twofold. [S4]* | 19 (20.2%) | [S1], [S2], [S4], [S13], [S14], [S15], [S20], [S32], [S43], [S53], [S56], [S62], [S64], [S69], [S72], [S73], [S84], [S86], [S93] |
| Improving the Adaptability of Agent System | *Our proposal for dynamic process generation aims to accommodate this variability, enabling a wider array of diverse instances to emerge and guide the development process accordingly. [S57]* | 18 (19.1%) | [S1], [S2], [S9], [S12], [S31], [S33], [S36], [S37], [S38], [S47], [S52], [S53], [S57], [S75], [S76], [S89], [S91], [S94] |
| Enhancing the Diversity of Generated Artifacts | *The parallel generation method with KMeans filtering helps improve the diversity than using parallel generation only, albeit not significantly. [S18]* | 12 (12.8%) | [S18], [S21], [S23], [S26], [S38], [S45], [S48], [S57], [S59], [S80], [S85], [S89] |
| Ensuring Software Security | *Risk-sensitivity: We guide the Planner with two pre-defined atomic actions to limit the action space of the planner: slicing and operation. [S22]* | 4 (4.3%) | [S22], [S77], [S82], [S83] |

three stages: "*Fault Comprehension*", "*Codebase Navigation*", and "*Fault Confirmation*". Then they implemented "*Document-Guided Search strategy*" and "*retained only the Top-N classes that possess relatively high method-level coverage*" to significantly reduce input scale to optimize context length [S4].

***Improving the Efficiency of Generated Software Artifacts*** (23, 24.5%) is also considered by lots of designers. The efficiency of producing software artifacts is a critical aspect that many designers of LLM-based MASs focus on, and the efficiency of production is a vital QA of MAS. Consequently, various designers use this as the design rationale to develop and implement LLM-based MASs. For example, Lee *et al.* proposed an LLM-based MAS called "*FixAgent*" under the guidance of *Hierarchical Coordination* to finish the bug fixing tasks using the least calculation resource [S2].

***Reducing the Difficulty of Task Resolution*** (19, 20.2%) is also a rationale considered by various designers. Reducing the difficulty in solving SE tasks can facilitate the design and implementation of MAS based on LLM, improve the quality of generated software artifacts, improve the efficiency of the production of software artifacts, and lower the costs associated with the completion of SE tasks. For instance, Qian *et al.* proposed an LLM-based MAS called "*ChatDev*", which "*uses a chat chain to divide each phase into smaller subtasks further*" to reduce the difficulty of task solving [S93].

***Improving the Adaptability of Agent System*** (18, 19.1%) is also a widely used rationale. Numerous designers of LLM-based MASs for SE tasks want to enable these systems to operate in diverse runtime environments, thereby placing strong emphasis on enhancing the adaptability of the MAS.

For instance, Xu *et al.* proposed an LLM-based MAS called "*Gentopia.AI*", which allows agents to be built in a hierarchical architecture and a non-hierarchical architecture, "*thereby extending their capabilities beyond single language models*" [S76].

**Enhancing the Diversity of Generated Software Artifacts** (12, 12.8%) is also considered by some designers. The diversity of software artifacts generated by LLM-based MASs is crucial. Enhancing the diversity of produced artifacts allows selection from a broader range, facilitating the identification of higher-quality artifacts and increasing the likelihood of finding optimal task solutions. For example, Ataei *et al.* proposed an LLM-based MAS called "*Elicitron*", which implements an automated and highly scalable "*process of creating and interviewing agents*", so that the MAS can "*identify a diverse set of user needs*" [S48].

**Ensuring Software Security** (4, 4.3%) also cannot be ignored. The security of LLM-based MASs for SE tasks is of paramount importance. Some designers of LLM-based MASs have incorporated specialized measures to ensure system security. For example, Khonee *et al.* proposed an LLM-based MAS called "*GoNoGo*", which "*guides the Planner with two pre-defined atomic actions to limit the action space of the planner*" to let the MAS be "*risk-sensitive*" [S22].

## 5 Discussions

### 5.1 Interpretations of the Results

In this section, we discuss and explain the relationship between the results of the four RQs: *SE Tasks*, *Quality Attributes*, *Design Patterns*, and *Design Rationale*.

*5.1.1 Mapping of SE Tasks and Quality Attributes.* Table 6 shows the mapping relationship between the *Quality Attribute* categories and *SE Task* categories in LLM-based MASs, using abbreviations to represent each category of *SE Tasks*. For example, "CG" represents *Code Generation*. The full names of all cause categories are provided in the note of Table 6.

From the results presented in Table 6, we can see that: *Functional Suitability* (89, 94.8%) is the QA considered by most of designers. Almost every LLM-based MAS has the design to maintain the *Functional Suitability* of the MAS. In the *Functional Suitability* category, *Functional Correctness* (86) is the most considered Sub-Quality Attribute (Sub-QA), which is paid most attention to by the designers of LLM-based MASs for *Code Generation* (43). As generated code directly determines whether a system behaves correctly, designers of LLM-based MASs for *Code Generation* prioritize *Functional Correctness* to ensure outputs conform to specifications, reduce debugging and integration costs. Besides, *Functional Correctness* is considered by designers of LLM-based MASs for every category of the *SE Task*, which indicates its highest importance. *Functional Completeness* is also considered by various designers, which is very important in systems for *Code Generation* (13). *Functional Appropriateness* is less considered in *Functional Suitability* category, but it cannot also be ignored. The predominance of *Functional Suitability*, especially *Functional Correctness*, reflects the imperative of designers to ensure that LLM-based MASs reliably produce semantically accurate and syntactically valid artifacts across all SE tasks.

Designers of LLM-based MASs for SE tasks also highly value *Performance Efficiency* (48, 51.1%). In *Performance Efficiency* category, both *Resource Utilization* and *Time Behaviour* are placed of great importance by designers of LLM-based MASs. *Resource Utilization* is highly valued by the designers of MASs for *Code Generation* (12). *Capability* is less considered in *Performance Efficiency* category. Given that LLM-based MASs have significant computational and latency costs, designers prioritize resource utilization and runtime behavior to ensure the performance of MASs. In addition, tasks such as *Code Generation* and *Fault Localization* are highly sensitive to response time: These

two SE tasks are critical steps in developing workflows, and high latency directly disrupts these processes. These tasks often require frequent model invocations (e.g., multiple completions, iterative localization and verification). When serving plenty of users or large codebases, resource efficiency becomes the key determinant of scalability and economic viability of an MAS. Although model capability remains important, in these engineering contexts the marginal gains from improving capability typically demand substantial additional computation and latency (e.g., larger models), so capability is often treated as a secondary priority.

*Maintainability* (47, 50.0%) is also a very significant QA considered by designers of LLM-based MASs for SE tasks. Moreover, *Modularity* is the most valued Sub-QA in *Maintainability* category, which is attached great importance in LLM-based MASs for *Code Generation* (20). However, *Analysability* and *Reusability* are both less considered in the design process for LLM-based MASs. The strong emphasis on *Maintainability* in LLM-based MAS design, especially *Modularity*, stems from the need to decompose complex agent workflows into interchangeable components. Implement agents as encapsulated, interchangeable modules with clearly defined interfaces, and assign each module a specific role and task. Such role-level modularity reduces coupling, constrains failure propagation, and enables independent development, testing, and replacement of agents. In addition, partitioning functionality into agent modules simplifies monitoring and analysis of individual components and produces components that are more readily reusable in other contexts. For these reasons, *Modularity* receives greater emphasis than *Analysability* and *Reusability*.

*Reliability* (34, 36.2%) is also considered by various designers of LLM-based MASs for SE tasks. *Fault Tolerance* is more valued than *Faultlessness* in the design of LLM-based MASs. *Fault Tolerance* is most valued in the design of systems for *Code Generation* (9). The emphasis on *reliability* arises because designers hope that agents based on LLM MASs must handle the inherent unpredictability and error proneness of generated artifacts to ensure robust workflows. As LLMs have exhibited intrinsic stochasticity, distributional sensitivity, and occasional hallucinations, making perfect correctness infeasible. Therefore designers emphasize redundancy, error detection, and graceful degradation to preserve system-level functionality when individual agents produce uncertain or incorrect outputs, which is emphasized in MASs for code generation. Designers prioritize iterative validation, automated testing, and fallback mechanisms that recover from erroneous outputs rather than relying on faultless generation.

*Flexibility* (24, 25.5 %) is also a QA considered by designers of LLM-based MASs. *Adaptability* is more considered by designers than *Scalability*. *Adaptability* is most valued QA in LLM-based MASs for *Code Generation* (8). Emphasis on *Flexibility* reflects the need for LLM-based MASs to adjust dynamically to evolving requirements and contexts from users. Designers prioritize *Adaptability* over *Scalability* because LLM-based systems must operate across evolving APIs, and shifting prompt distributions, requiring agents to rapidly adjust behaviors, representations, and tool use to provide correct function. This preference is most evident in LLM-based MASs for *Code Generation*, which have frequent changes in programming languages, libraries, and coding conventions necessitate continual adaptation.

Many designers of LLM-based MASs for SE tasks also consider *Compatibility* (10, 10.6%) of their systems for SE tasks. *Interoperability* is more considered by designers than *Co-Existence*. *Interoperability* is more valued in the systems for *Code Generation* (6). Emphasis on *Compatibility* stems from the necessity for LLM-based MAS to integrate seamlessly with existing development environments. *Interoperability* enables semantically aligned communication, capability composition, and end-to-end verification among agents, so this attribute is more valued than *Co-Existence*. There are interoperable toolchains, common intermediate representations (e.g., ASTs), and consistent formatting and typing conventions in systems for *Code Generation*, so it is important to ensure *Compatibility*.

*Security* (10, 10.6%) is likewise considered by some designers of LLM-based MASs for SE tasks and shares the same proportion as *Compatibility*. *Confidentiality* is the most considered Sub-QA in *Security* category. *Confidentiality* is most valued in the systems for *Code Generation* (2), *Fault Localization* (2), and *End-to-End Software Development* (2). Emphasis on *Security* reflects growing recognition of the risks associated with sensitive data exposure during automated code synthesis and analysis from designers of LLM-based MASs. Designers prioritize *Confidentiality* within the security category because LLM-based systems routinely process proprietary code, sensitive datasets, and user secrets. In addition, workflows expose intellectual property, debugging traces, test fixtures, and credentials its leakage can cause competitive loss, security incidents, and regulatory penalties. Therefore, it is important to ensure *Security* in systems for *Code Generation*, *Fault Localization*, and *End-to-End Software Development*.

Some designers of LLM-based MASs for SE tasks also want to have *Interaction Capability* (9, 9.6%) in their systems. *Operability* shares the same proportion as *User Engagement* in *Interaction Capability* category. *Operability* is most valued in systems for *Code Generation* (4), and *User Engagement* is also most considered in systems for *Code Generation* (3). Emphasis on *Interaction Capability* stems from the inherently iterative nature of LLM workflows, where clear operability and active user engagement are essential for soliciting feedback and refining agent outputs. Effective human–system interaction requires both user-facing qualities (e.g., clarity and intuitive control) and system-facing qualities (e.g., observability and controllability), balancing these qualities ensures agents are approachable to users while remaining manageable and maintainable in real-world deployments. Therefore, *User Engagement* and *Operability* are both important.

Table 6. Mapping between Quality Attributes (vertical) and SE Tasks (horizontal)

| | | CG | FL | ETESM | PR | ETESD | CR | ST | RE | CT | SD |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Functional Suitability** | Functional Correctness | 43 | 8 | 7 | 8 | 6 | 4 | 6 | 2 | 1 | 1 |
| | Functional Completeness | 13 | 0 | 0 | 0 | 2 | 3 | 4 | 2 | 0 | 0 |
| | Functional Appropriateness | 6 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 |
| **Performance Efficiency** | Resource Utilization | 12 | 5 | 3 | 1 | 1 | 2 | 1 | 2 | 0 | 0 |
| | Time Behaviour | 8 | 2 | 4 | 2 | 4 | 2 | 3 | 0 | 1 | 1 |
| | Capability | 3 | 1 | 1 | 1 | 1 | 1 | 1 | 2 | 0 | 0 |
| **Maintainability** | Modularity | 20 | 4 | 5 | 3 | 3 | 5 | 3 | 2 | 1 | 0 |
| | Analysability | 3 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| | Reusability | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **Reliability** | Fault Tolerance | 9 | 0 | 1 | 2 | 4 | 3 | 1 | 1 | 0 | 0 |
| | Faultlessness | 10 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| **Flexibility** | Adaptability | 8 | 1 | 2 | 0 | 2 | 1 | 2 | 0 | 0 | 0 |
| | Scalability | 7 | 0 | 1 | 3 | 1 | 0 | 0 | 0 | 0 | 0 |
| **Compatibility** | Interoperability | 6 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| | Co-Existence | 4 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| **Security** | Confidentiality | 2 | 2 | 1 | 0 | 2 | 1 | 1 | 0 | 0 | 0 |
| | Integrity | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **Interaction Capability** | Operability | 4 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| | User Engagement | 3 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |

**Full names of each cause type:** CG: *Code Generation*; FL: *Fault Localization*; ETESM: *End-to-End Software Maintenance*; PR: *Program Repair*; ETESD: *End-to-End Software Development*; CR: *Code Review*; ST: *Software Testing*; RE: *Requirements Engineering*; CT: *Code Translation*; SD: *Software Deployment*.

*5.1.2 Mapping of SE Tasks, Quality Attributes, Sub-Quality Attributes, Design Patterns, and Design Rationale.* Figure 4 shows the mapping relationship between the *SE Task* categories, *Quality Attributes* categories, *Sub-Quality Attributes* categories, *Design Pattern* categories, and *Design Rationale* categories in LLM-based MASs for SE tasks.
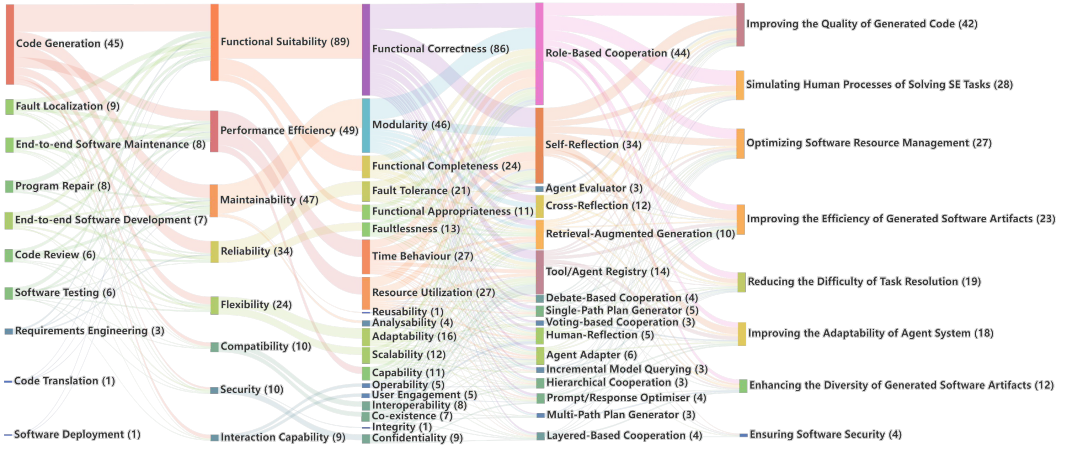
Fig. 4. Mapping between SE Tasks, Quality Attributes, Sub-Quality Attributes, Design Patterns, and Design Rationale

From the results shown in Figure 4, we can see that most of the designers of LLM-based MASs use *Role-Based Cooperation* (18) and *Self-Reflection* (18) to construct the systems for *Code Generation*, and the most used guidance is *Improving the Quality of Generated Code* (25). *Code Generation* critically depends on clear role specialization for modular task decomposition and self-reflective feedback loops, all aimed at maximizing the correctness and maintainability of the generated code. Therefore, a lot of designers use *Role-Based Cooperation* and *Self-Reflection* to implement their LLM-based MASs, thereby improving the quality of generated code.

Many designers of LLM-based MASs use *Role-Based Cooperation* (6) to build the systems for *Fault Localization*, and the most used guidance is *Optimizing Software Resource Management* (5). It is because assigning specialized roles enables parallel exploration of fault candidates while resource-management optimizations minimize the overhead of coordinating multiple diagnostic agents, various designers like to use *Role-Based Cooperation* for *Fault Localization*.

Designers of LLM-based MASs for *End-to-End Software Maintenance* are likely to use *Role-Based Cooperation* (6) to construct their systems, and the most used guidance is *Improving the Quality of Generated Code* (3), *Simulating Human Processes of Solving SE Tasks* (3), *Optimizing Software Resource Management* (3), and *Improving the Efficiency of Generating Software Artifacts* (3). *End-to-End Software Maintenance* demands a clear division of specialized responsibilities, assigning specialized roles (e.g., fault localizer, patch generator, tester, and integrator) modularizes complex diagnostic and repair workflows, reduces interference between competing model behaviors, improves traceability and accountability, and enables parallelism and targeted fine-tuning. Designers emphasize code quality, human problem-solving processes, resource management, and efficiency, because these priorities collectively enhance interpretability and traceability, control computational cost, and accelerate the reliable delivery of generated artifacts.

Many designers of LLM-based MASs use *Role-Based Cooperation* (2), *Cross-Reflection* (2), and *RAG* (2) to construct the systems for *Program Repair*, and the most used guidance is *Improving the Quality of Generated Code* (3). Because role specialization decomposes the repair workflow into modular responsibilities (e.g., synthesizer, tester, validator), cross-reflection enables iterative peer critique that exposes and corrects mistakes, and RAG grounds generations in external code and test

artifacts to improve factuality and relevant recall, the aforementioned patterns collectively support semantic correctness, syntactic validity, and low regression risk.

Designers of LLM-based MASs for *End-to-End Software Development* are willing to use *Role-Based Cooperation* (6) to build their systems, and the most used guidance is *Improving the Quality of Generated Code* (6). By creating a well-defined division of labor among agents and prioritizing generated code quality, *Role-Based Cooperation* promotes *End-to-End Software Development* pipelines that yield correct, maintainable artifacts.

Many designers of LLM-based MASs use *Role-Based Cooperation* (2) to build the systems for *Code Review*, and the most used guidance is *Simulating Human Processes of Solving SE Tasks* (4). *Code Review* inherently involves distinct reviewer roles and benefits from modeling human critique workflows, making *Role-Based Cooperation* combined with human-process simulation an effective design strategy.

Designers of LLM-based MASs for *Software Testing* are likely to use *Self-Reflection* (3) to construct their systems, and the most used guidance is *Improving the Quality of Generated Code* (3). *Software Testing* tasks can benefit from the ability of agents to introspect and refine outputs, while emphasizing code quality, adaptability, and diversity, aligning with the goals of generating robust and comprehensive test artifacts.

Designers of LLM-based MASs for *Requirements Engineering* are willing to use *Role-Based Cooperation* (3) to build their systems, and the most used guidance is *Optimizing Software Resource Management* (2), and *Enhancing the Diversity of Model Generation* (2). Because assigning specialized roles allows coordinated, parallel handling of elicitation, negotiation, and validation tasks, which reduces redundant computation and resource contention, while role specialization and varied agent perspectives increase the breadth of candidate solutions and reduce generation collapse.

Designers of LLM-based MASs use *Role-Based Cooperation* (1) to construct the systems for *Code Translation*, and the most used guidance is *Reducing the Difficulty of Task Resolution* (1). *Code Translation* is a cognitively intensive task involving complex syntax and semantics, so that assigning specialized roles helps decompose the task and reduce its resolution difficulty.

Designers of LLM-based MASs use *Self-Reflection* (1) and *Agent Adapter* (1) to build the systems for *Software Deployment*, and the most used guidance is *Ensuring Software Security* (1). *Software Deployment* requires agents to autonomously evaluate their outputs and adapt to diverse deployment environments, so that designers can use *Self-Reflection* and *Agent Adapter* to design their LLM-based MASs. Besides, the reflection mechanism lets agents detect and correct deployment-time faults autonomously, while adapters provide modular, environment-specific integration, and together these practices prioritize ensuring software security to mitigate the enlarged attack surface and operational risks introduced by adaptive, automated agents.

## 5.2 Implications

In this section, we provide the implications of this study based on the results in Section 4 and the relationships between the results of the four RQs in Section 5.1, intended to serve as practical guidance for the design and implementation of LLM-based MASs for SE tasks.

> **Implication 1**. Designers need to place greater emphasis on the quality of generated software artifacts when designing LLM-based MASs.

Designers of LLM-based MASs should place a significant emphasis on *Functional Correctness* (91.5%), which is the *Quality Attribute* most emphasized by designers of LLM-based MASs. *Functional Correctness* directly affects whether generated artifacts meet specification requirements and perform the intended tasks reliably. Therefore, the emphasis on artifact quality should be reflected in concrete

design choices. For example, designers could adopt rigorous specification and validation practices, including formal or semi-formal specifications where feasible, extensive unit and integration tests driven by both automated test generation and manually curated test cases, and continuous verification in realistic execution environments.

Besides, *Improving the Quality of Generated Code* (44.7%) is the most considered *Design Rationale*. Currently, *Code Generation* is the SE task most frequently addressed by LLM-based MASs. Therefore, designers place special emphasis on the quality of the generated code. However, more and more SE tasks are solved by specially constructed MASs, quality assurance for other types of generated software artifacts cannot be neglected. The quality of the generated artifacts determines whether those artifacts satisfy functional requirements, integrate cleanly with existing systems, and remain maintainable under realistic operating conditions. For example, Nguyen *et al.* [S5] partition the development workflow into specialized roles (i.e., PM, SM, Developer, Senior Developer, and Tester) and employ a three-step static code review procedure to enhance functional correctness, while leveraging a dynamic code dependency graph (DCGG) to maintain inter-file dependencies and thereby improve overall code quality.

---

**Implication 2**. Designers could use *Role-Based Cooperation* to improve the *Maintainability* of LLM-based MASs.

---

In Figure 4, *Maintainability* of LLM-based MASs primarily stems from *Modularity*, which is largely derived from *Role-Based Cooperation*. Therefore, designers of LLM-based MASs for SE tasks can employ *Role-Based Cooperation* within the systems to enhance maintainability by decoupling concerns, encapsulating functionality into specialized agents, and enabling targeted updates, testing, and evolution of individual system components. Under the guidance of *Role-Based Cooperation*, each agent encapsulates its own logic, prompt strategies, and tool interfaces. The clear separation of roles enables agents to be developed, debugged, and tested independently, thereby enhancing modularity. Moreover, the modular role structure allows individual agents to be replaced or upgraded without altering the overall system architecture. Additionally, system behaviors are distributed across distinct roles, which facilitates tracking of issues to specific agents. Therefore, *Role-Based Cooperation* mechanism not only improves the structural clarity of LLM-based MASs but also significantly enhances system maintainability. For example, Qin *et al.* [S4] proposed an LLM-based MAS called *AGENTFL*, which incorporates four different agents (i.e., Test Code Reviewer, Source Code Reviewer, Software Architect, and Software Test Engineer) to achieve high modularity and maintainability, enabling each agent to specialize, evolve, and interact in a structured and extensible way.

---

**Implication 3**. Designers can draw inspiration from human-centered software development activities to design their systems when developing LLM-based MASs.

---

Our analysis shows that *Simulating Human Processes of Solving SE Tasks* (29.8%) is a common rationale adopted by various designers of LLM-based MASs, which indicates that designers could leverage human-centered software development methodologies when building LLM-based MASs for SE tasks. Drawing on human-centered software development practices, designers can adopt established principles of task decomposition, role specialization, coordination, and verification to enhance interpretability, verifiability, and robustness of an MAS. Designers could study human cooperative development patterns and human cognitive processes (e.g., task decomposition, distributed attention, working memory limits, and error-detection strategies) and translate those insights into explicit role-specialization, hierarchical task decomposition, and lightweight coordination protocols

that include built-in verification checkpoints. Certainly, there are numerous challenges, such as the mismatch between human assumptions and agent capabilities, ambiguities in communication protocols, and the scarcity of methods to validate the results and quality of interactions among agents. However, these issues can be mitigated by adopting iterative "human-in-the-loop" evaluation approaches, architectures inspired by human cognitive processes, etc. For example, Lee *et al.* [S2] designed their LLM-based MAS for software debugging using a methodology named *Hierarchical Cooperation*. They structured multiple specialized agents into a tiered coordination framework, where each level addresses faults of increasing complexity through progressively sophisticated strategies like the human cognitive model of debugging. By emulating mature human SE activities, the complexity of designing LLM-based MASs can be reduced. Agents are organized into a collaborative development team, where complex tasks are decomposed into specialized subtasks and assigned to dedicated agents, each of which performs its role, exploits its strengths, and together drives the software development process just like human beings.

> **Implication 4**. Designers are beginning to leverage LLM-based MASs for the entire lifecycle support.

The result that *End-to-end Maintenance* (8.4%) and *End-to-End Software Development* (7.4%) are resolved by various, especially constructed LLM-based MASs suggests that designers are increasingly adopting LLM-based MASs to support the entire software lifecycle. A unified MAS retains contextual knowledge from requirements through to final software, ensuring that design decisions, coding conventions, and test criteria remain coherent across all phases. In our dataset, the first paper proposing an LLM-based MAS for *End-to-End Software Development* was published in July 2023, with only four such papers published throughout that year. However, by October 2024, the number of related works had risen to seven. Besides, all papers focusing on LLM-based MASs for *End-to-end Maintenance* were published exclusively in 2024. The relative paucity of papers on end-to-end development and maintenance stems from lots of reasons. First, there is a lack of comprehensive evaluation criteria that can meaningfully assess system behavior across requirements engineering, software architecture design, implementation, testing and maintenance. Without such unified metrics, empirical comparison and validation are difficult. Second, there is an absence of complete ground-truth guidance for each development stage. Although several benchmarks exist for code-generation (e.g., HumanEval and MBPP), and recent efforts such as DevBench have begun to target multiple stages of the software lifecycle, there remains no widely accepted ground-truth benchmark that comprehensively evaluates multi-agent systems across the full software development lifecycle. Third, achieving alignment across stages is intrinsically hard: models or tools optimized for one task often produce outputs that are incompatible with the assumptions or interfaces of downstream stages, so end-to-end coordination requires solving difficult cross-stage alignment problems. However, despite these obstacles, the growing research and tool support for individual SE tasks is steadily producing the theoretical foundations and technical components needed for end-to-end systems, and designers are increasingly willing to build LLM-based MASs for full-lifecycle software development and maintenance. For example, Sami *et al.* [S6] introduce a unified platform that deploys specialized AI agents tailored to specific tasks, including the generation of user stories, prioritization of requirements, creation of UML diagrams, code generation, and automated testing to fill the absence of a cohesive platform capable of delivering consistent and optimal results across all phases of the software development lifecycle.

**Implication 5**. The rationale for resource-oriented and efficiency-oriented design reflects the considerations of designers to minimize temporal and spatial costs when completing relevant SE tasks.

The significant emphasis on *Optimizing Software Resource Management* (28.7%) and *Improving the Efficiency of Software Artifact Generation* (24.5%) underscores that designers cannot ignore the considerations of temporal and spatial costs when guiding MAS design. Besides, *Time Behaviour* (27) and *Resource Utilization* (27) are considered by many designers of LLM-based MASs. These two QAs directly determine the practicality and effectiveness of an MAS. Poor time behaviour undermines user experience and can violate real-time requirements for interactive and safety-critical tasks, while inefficient resource usage raises operational costs, limits scalability, and may preclude deployment on constrained hardware. Therefore, designers of LLM-based MASs think of various methods to reduce temporal and spatial costs. For example, Chen *et al.* [S14]. introduces a novel multi-agent framework named *CODER*, which is augmented with a task graph data structure to systematically resolve GitHub issues within software repositories. By encoding each debug workflow as a JSON-formatted task graph, *CODER* eliminates redundant plan synthesis. All agents in the MAS follow a strictly executable plan, avoiding iterative LLM prompts and context reloading that incur latency and API costs.

**Implication 6**. Designers could focus on the performance efficiency of LLM-based MASs, as well as consider their maintainability.

Our data analysis results show that 51.1% of LLM-based MASs explicitly consider *Performance Efficiency* during design, while 50.0% of designers reported prioritizing *Maintainability*. Besides, there are 21 LLM-based MASs (22.3%) that are built under the consideration of both *Performance Efficiency* and *Maintainability*, which indicates that system designers could concurrently ensure both immediate operational efficiency and sustained maintainability. Treating performance and maintainability as concurrent design objectives for LLM-based MASs promotes sustained operational efficiency. Many designers have made their efforts for this aim. Designers can trade off competing QAs in LLM-based MASs by making role assignment and coordination strategies first-class design decisions. For example, assigning specialized agents with narrow, well-defined responsibilities, enables the MAS to satisfy performance-oriented metrics such as time behavior and resource utilization for routine cases while reserving more expensive, higher-assurance procedures for difficult cases (e.g., FixAgent proposed by Lee *et al.* [S2]). Hence, employing well-structured role assignments together with robust coordination protocols enables designers to negotiate conflicts among QAs and to systematically determine trade-off policies. In addition, Shen *et al.* [22] presented an empirical study on LLM-based MASs for software development and gave an optimization methodology for these systems using textual feedback. They introduced a role-based cooperation workflow composed of a recruitment agent (for static role assignment), a solving agent (developer), and a review agent (test engineer), executed in an iterative optimization cycle of recruitment, solving, review, and subsequent solving. Additionally, time efficiency is improved through One-Pass Prompting and Offline Feedback Caching, while space efficiency is optimized via Prompt Compression and Trajectory Truncation.

## 6 Threats on Validity

In this section, we outline the potential threats to the validity of our study. We identify the threats encountered during the research and clarify the measures we employed to mitigate them. Internal

validity is not discussed, as our study does not involve any experimental manipulation of variables and thus does not support causal inference.

**Construct validity**: Since data collection and data extraction in our study were conducted manually, there is a potential risk of individual bias in the data extraction results. To mitigate this risk, a pilot data extraction was conducted before the formal data extraction, which partially alleviated the threats to the construct validity of the study. Furthermore, following both the pilot and formal data extraction, the first author engaged discussions with the second and third authors to ensure that all of the data items could be extracted from our dataset and the data extraction results were aligned with four RQs.

**External validity**: In this study, external validity primarily concerns the selection of data sources. To ensure the credibility of our data, our data collection is based on two recent literature surveys on LLM-based agent systems for SE tasks by Liu *et al.* [17] and Wang *et al.* [26]. To further enhance the comprehensiveness of the dataset, we additionally included arXiv [25] as a data source, which is an open-access preprint platform maintained by Cornell University and is widely recognized within the academic community. Besides, arXiv features a dedicated "Software Engineering" category, which facilitates the identification of papers relevant to our RQs.

**Reliability**: To mitigate potential uncertainties associated with the adopted research methodology, we implemented several measures to enhance the reliability of the study. Throughout the processes of data extraction and data analysis, extensive discussions were held among the first, second, and third authors to resolve any internal inconsistencies and to ensure the consistency and accuracy of the results. Moreover, we have made our dataset [3] publicly available to enable other researchers to replicate the study and validate our findings.

## 7   Conclusions and Future Work

In this study, we focused on the SE tasks addressed by the specially designed LLM-based MASs, as well as the *quality attributes* considered by the designers of LLM-based MASs to address the SE tasks, the *design patterns* employed to build LLM-based MASs for SE tasks, and the *design rationale* supporting the construction of LLM-based MASs to facilitate SE tasks. We collected 94 papers that met our criteria from two recent surveys on LLM-based agent systems for SE tasks by Liu *et al.* [17] and Wang *et al.* [26], and the SE category of arXiv [25]. The study results show that: *Code Generation* is the most common SE task addressed by LLM-based MASs, *Functional Suitability* is the QA which is mostly considered by designers of LLM-based MASs for SE tasks, *Role-Based Cooperation* is the most frequently used design pattern to develop LLM-based MASs to support SE tasks, and *Improving the Quality of Generated Code* is the most common rationale behind the design of LLM-based MASs for SE tasks.

Based on our study results, we provide implications for designing LLM-based MASs for SE tasks. For example, designers could use *Role-Based Cooperation* to improve the *Maintainability* of LLM-based MASs. In addition, designers can draw inspiration from human-centered software development activities to design their systems when developing LLM-based MASs. Moreover, designers are beginning to leverage LLM-based MASs for supporting the entire software development lifecycle. The rationale behind resource-oriented and efficiency-oriented designs reflects designers' intent to minimize temporal and spatial costs when completing relevant SE tasks.

## Data Availability

The dataset of this work has been made available at [3].

## Acknowledgments

## References

[1] Chetan Arora, John Grundy, and Mohamed Abdelrazek. 2023. Advancing Requirements Engineering through Generative AI: Assessing the Role of LLMs. *arXiv preprint arXiv:2310.13976* (2023).

[2] Islem Bouzenia and Michael Pradel. 2025. Understanding Software Engineering Agents: A Study of Thought-Action-Result Trajectories. *arXiv preprint arXiv:2506.18824* (2025).

[3] Yangxiao Cai, Ruiyin Li, Peng Liang, Mojtaba Shahin, and Zengyang Li. 2025. *Dataset of the Paper "Designing LLM-based Multi-Agent Systems for Software Engineering Tasks: Quality Attributes, Design Patterns and Rationale"*. https://github.com/Caiyangxiao/MASDesign.

[4] Mert Cemri, Melissa Z. Pan, Shuyi Yang, Lakshya A. Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, Matei Zaharia, Joseph E. Gonzalez, and Ion Stoica. 2025. Why Do Multi-Agent LLM Systems Fail? *arXiv preprint arXiv:2503.13657* (2025).

[5] Shuaihang Chen, Yuanxing Liu, Wei Han, Weinan Zhang, and Ting Liu. 2024. A Survey on LLM-based Multi-Agent System: Recent Advances and New Frontiers in Application. *arXiv preprint arXiv:2412.17481* (2024).

[6] Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2024. AgentVerse: Facilitating Multi-Agent Collaboration and Exploring Emergent Behaviors. In *Proceedings of the 12th International Conference on Learning Representations (ICLR)*. ICLR, 1–43.

[7] Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. 2024. Self-Collaboration Code Generation via ChatGPT. *ACM Transactions on Software Engineering and Methodology* 33, 7 (2024), 1–38.

[8] Ali Dorri, Salil S. Kanhere, and Raja Jurdak. 2018. Multi-Agent Systems: A Survey. *IEEE Access* 6 (2018), 28573–28593.

[9] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. [n. d.]. Large Language Model Based Multi-agents: A Survey of Progress and Challenges. In *Proceedings of the 33rd International Joint Conference on Artificial Intelligence (IJCAI)*. IJCAI, 8048–8057.

[10] Junda He, Christoph Treude, and David Lo. 2025. LLM-Based Multi-Agent Systems for Software Engineering: Literature Review, Vision and the Road Ahead. *ACM Transactions on Software Engineering and Methodology* 34, 5 (2025), Article No.: 124.

[11] Sirui Hong, Mingchen Zhuge, Jiaqi 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. 2023. MetaGPT: Meta Programming for A Multi-Agent Collaborative Framework. *arXiv preprint arXiv:2308.00352* (2023).

[12] Jie Huang and Kevin Chen-Chuan Chang. 2023. Towards Reasoning in Large Language Models: A Survey. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL)*. ACL, 1049–1065.

[13] International Organization for Standardization and International Electrotechnical Commission. 2023. *ISO/IEC 25010:2023 Systems and software engineering – Systems and software Quality Requirements and Evaluation (SQuaRE) – System and software quality models*. Technical Report ISO/IEC 25010:2023. ISO/IEC, Geneva, Switzerland.

[14] Dongming Jin, Zhi Jin, Xiaohong Chen, and Chunhui Wang. 2024. MARE: Multi-Agents Collaboration Framework for Requirements Engineering. *arXiv preprint arXiv:2405.03256* (2024).

[15] R. Kazman, M. Klein, M. Barbacci, T. Longstaff, H. Lipson, and J. Carriere. 1998. The architecture tradeoff analysis method. In *Proceedings of 4th IEEE International Conference on Engineering of Complex Computer Systems (ICECCS)*. 68–78.

[16] Ruiyin Li, Yiran Zhang, Xiyu Zhou, Peng Liang, Weisong Sun, Jifeng Xuan, Zhi Jin, and Yang Liu. 2025. MAAD: Automate Software Architecture Design through Knowledge-Driven Multi-Agent Collaboration. *arXiv preprint arXiv:2507.21382* (2025).

[17] Junwei Liu, Kaixin Wang, Yixuan Chen, Xin Peng, Zhenpeng Chen, Lingming Zhang, and Yiling Lou. 2024. Large Language Model-Based Agents for Software Engineering: A Survey. *arXiv preprint arXiv:2409.02977* (2024).

[18] Yue Liu, Sin Kit Lo, Qinghua Lu, Liming Zhu, Dehai Zhao, Xiwei Xu, Stefan Harrer, and Jon Whittle. 2025. Agent Design Pattern Catalogue: A Collection of Architectural Patterns for Foundation Model based Agents. *Journal of Systems and Software* 220 (2025), 112278.

[19] Ruwei Pan, Hongyu Zhang, and Chao Liu. 2025. CodeCoR: An LLM-Based Self-Reflective Multi-Agent Framework for Code Generation. *arXiv preprint arXiv:2501.07811* (2025).

[20] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. ChatDev: Communicative Agents for Software

Development. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*. ACL, 15174–15186.

[21] Anjana Sarkar and Soumyendu Sarkar. 2025. Survey of LLM Agent Communication with MCP: A Software Design Pattern Centric Review. *arXiv preprint arXiv:2506.05364* (2025).

[22] Ming Shen, Raphael Shu, Anurag Pratik, James Gung, Yubin Ge, Monica Sunkara, and Yi Zhang. 2025. Optimizing LLM-Based Multi-Agent System with Textual Feedback: A Case Study on Software Development. *arXiv preprint arXiv:2505.16086* (2025).

[23] Wentao Shi, Xiangnan He, Yang Zhang, Chongming Gao, Xinyue Li, Jizhi Zhang, Qifan Wang, and Fuli Feng. 2024. Large Language Models are Learnable Planners for Long-Term Recommendation. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*. ACM.

[24] Klaas-Jan Stol, Paul Ralph, and Brian Fitzgerald. 2016. Grounded theory in software engineering research: a critical review and guidelines. In *Proceedings of the 38th International Conference on Software Engineering (ICSE)*. ACM, 120–131.

[25] Cornell University. 2025. *arXiv*. https://arxiv.org/list/cs.SE/recent.

[26] Yanlin Wang, Wanjun Zhong, Yanxian Huang, Ensheng Shi, Min Yang, Jiachi Chen, Hui Li, Yuchi Ma, Qianxiang Wang, and Zibin Zheng. 2024. Agents in Software Engineering: Survey, Landscape, and Vision. *Automated Software Engineering* 32 (2024), Article No.: 70.

[27] Bingyu Yan, Zhibo Zhou, Litian Zhang, Lian Zhang, Ziyi Zhou, Dezhuang Miao, Zhoujun Li, Chaozhuo Li, and Xiaoming Zhang. 2025. Beyond Self-Talk: A Communication-Centric Survey of LLM-Based Multi-Agent Systems. *arXiv preprint arXiv:2502.14321* (2025).

[28] Miao Yu, Fanci Meng, Xinyun Zhou, Shilong Wang, Junyuan Mao, Linsey Pang, Tianlong Chen, Kun Wang, Xinfeng Li, Yongfeng Zhang, Bo An, and Qingsong Wen. 2025. A Survey on Trustworthy LLM Agents: Threats and Countermeasures. *arXiv preprint arXiv:2503.09648* (2025).

[29] Zhiqiang Yuan, Weitong Chen, Hanlin Wang, Kai Yu, Xin Peng, and Yiling Lou. 2024. TRANSAGENT: An LLM-Based Multi-Agent System for Code Translation. *arXiv preprint arXiv:2409.19894* (2024).

## A  Included Studies

Table 7. List of the included studies in this work

| ID | Author, Publication Title, and Venue |
|---|---|
| [S1] | Sarah Fakhoury, Markus Kuppe, Shuvendu K. Lahiri, Tahina Ramananandro, Nikhil Swamy. **3DGen: AI-Assisted Generation of Provably Correct Binary Format Parsers**. arXiv preprint arXiv:2404.10362 |
| [S2] | Cheryl Lee, Chunqiu Steven Xia, Longji Yang, Jen-tse Huang, Zhouruixin Zhu, Lingming Zhang, Michael R. Lyu. **FixAgent: Hierarchical Multi-Agent Framework for Unified Software Debugging** . arXiv preprint arXiv:2404.17153 |
| [S3] | Lyuye Zhang, Kaixuan Li, Kairan Sun, Daoyuan Wu, Ye Liu, Haoye Tian, Yang Liu. **ACFIX: Guiding LLMs with Mined Common RBAC Practices for Context-Aware Repair of Access Control Vulnerabilities in Smart Contracts**. arXiv preprint arXiv:2403.06838 |
| [S4] | Yihao Qin, Shangwen Wang, Yiling Lou, Jinhao Dong, Kaixin Wang, Xiaoling Li, Xiaoguang Mao. **AGENTFL: Scaling LLM-based Fault Localization to Project-Level Context**. arXiv preprint arXiv:2403.16362 |
| [S5] | Minh Huynh Nguyen, Thang Phan Chau, Phong X. Nguyen, Nghi D. Q. Bui. **AgileCoder: Dynamic Collaborative Agents for Software Development based on Agile Methodology**. arXiv preprint arXiv:2406.11912 |
| [S6] | Zeeshan Rasheed, Malik Abdul Sami, Muhammad Waseem, Kai-Kristian Kemell, Xiaofeng Wang, Anh Nguyen, Kari Systä, Pekka Abrahamsson. **AI-powered Code Review with LLMs: Early Results**. arXiv preprint arXiv:2404.18496 |
| [S7] | Bin Lei, Yuchen Li, Qiuwu Chen. **AutoCoder: Enhancing Code Large Language Model with AIEV-INSTRUCT**. arXiv preprint arXiv:2405.14906 |
| [S8] | Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, Abhik Roychoudhury. **AutoCodeRover: Autonomous Program Improvement**. In Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA), ACM, 1592 - 1604. |
| [S9] | Michele Tufano, Anisha Agarwal, Jinu Jang, Roshanak Zilouchian Moghaddam, Neel Sundaresan. **AutoDev: Automated AI-Driven Development**. arXiv preprint arXiv:2403.08299. |
| [S10] | Chenyuan Yang, Xuheng Li, Md Rakib Hossain Misu, Jianan Yao, Weidong Cui, Yeyun Gong, Chris Hawblitzel, Shuvendu Lahiri, Jacob R. Lorch, Shuai Lu, Fan Yang, Ziqiao Zhou, Shan Lu. **AutoVerus: Automated Proof Generation for Rust Code**. arXiv preprint arXiv:2409.13082. |
| [S11] | Sanjiban Choudhury, Paloma Sodhi. *Better than Your Teacher: LLM Agents that learn from Privileged AI Feedback*. arXiv preprint arXiv:2410.05434 |
| [S12] | Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, Zhi Jin. **CodeAgent: Enhancing Code Generation with Tool-Integrated Agent Systems for Real-World Repo-level Coding Challenges**. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL), ACL, 13643 - 13658 |
| [S13] | Zeeshan Rasheed, Malik Abdul Sami, Kai-Kristian Kemell, Muhammad Waseem, Mika Saari, Kari Systä, Pekka Abrahamsson. **CodePori: Large Scale Model for Autonomous Software Development by Using Multi-Agents**. arXiv preprint arXiv:2402.01411 |
| [S14] | Dong Chen, Shaoxin Lin, Muhan Zeng, Daoguang Zan, Jian-Gang Wang, Anton Cheshkov, Jun Sun, Hao Yu, Guoliang Dong, Artem Aliev, Jie Wang, Xiao Cheng, Guangtai Liang, Yuchi Ma, Pan Bian, Tao Xie, Qianxiang Wang. **CODER: ISSUE RESOLVING WITH MULTI-AGENT AND TASK GRAPHS**. arXiv preprint arXiv:2406.01304 |

| ID | Author, Publication Title, and Venue |
|---|---|
| [S15] | 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, Lizhen Cui. **CodeS: Natural Language to Code Repository via Multi-Layer Sketch**. arXiv preprint arXiv:2403.16443 |
| [S16] | Wei Ma, Daoyuan Wu, Yuqiang Sun, Tianwen Wang, Shangqing Liu, Jian Zhang, Yue Xue, Yang Liu. **Combining Fine-tuning and LLM-based Agents for Intuitive Smart Contract Auditing with Justifications**. arXiv preprint arXiv:2403.16073 |
| [S17] | Kexun Zhang, Weiran Yao, Zuxin Liu, Yihao Feng, Zhiwei Liu, Rithesh Murthy, Tian Lan, Lei Li, Renze Lou, Jiacheng Xu, Bo Pang, Yingbo Zhou, Shelby Heinecke, Silvio Savarese, Huan Wang, Caiming Xiong. **DIVERSITY EMPOWERS INTELLIGENCE:INTEGRAT-ING EXPERTISE OF SOFTWARE ENGINEERING AGENTS**. arXiv preprint arXiv:2408.07060 |
| [S18] | Mohammadmehdi Ataei, Hyunmin Cheong, Daniele Grandi, Ye Wang, Nigel Morris, Alexander Tessier. **Elicitron: An LLM Agent-Based Simulation Framework for Design Requirements Elicitation**. arXiv preprint arXiv:2404.16045 |
| [S19] | Md Nakhla Rafi, Dong Jae Kim, Tse-Hsun Chen, Shaowei Wang. **Enhancing Fault Localization Through Ordered Code Analysis with LLM Agents and Self-Reflection**. arXiv preprint arXiv:2409.13642 |
| [S20] | Simiao Zhang, Jiaping Wang, Guoliang Dong, Jun Sun, Yueling Zhang, Geguang Pu. **Experimenting a New Programming Practice with LLMs**. arXiv preprint arXiv:2401.01062 |
| [S21] | Malik Abdul Sami, Muhammad Waseem, Zeeshan Rasheed, Mika Saari, Kari Systä, Pekka Abrahamsson. **Experimenting with Multi-Agent Software Development: Towards a Unified Platform**. arXiv preprint arXiv:2406.05381 |
| [S22] | Arsham Gholamzadeh Khoee, Yinan Yu, Robert Feldt, Andris Freimanis, Patrick Andersson Rhodin, Dhasarathy Parthasarathy. *GoNoGo: An Efficient LLM-based Multi-Agent System for Streamlining Automotive Software Release Decision-Making*. arXiv preprint arXiv:2408.09785 |
| [S23] | Yingwei Ma, Qingping Yang, Rongyu Cao, Binhua Li, Fei Huang, Yongbin Li. **How to Understand Whole Software Repository?**. arXiv preprint arXiv:2406.01422 |
| [S24] | Chen Qian, Jiahao Li, Yufan Dang, Wei Liu, YiFei Wang, Zihao Xie, Weize Chen, Cheng Yang, Yingli Zhang, Zhiyuan Liu, Maosong Sun. **Iterative Experience Refinement of Software-Developing Agents**. arXiv preprint arXiv:2405.04219 |
| [S25] | Richard Fang, Rohan Bindu, Akul Gupta, Daniel Kang. **LLM Agents can Autonomously Exploit One-day Vulnerabilities**. arXiv preprint arXiv:2404.08144 |
| [S26] | Zhiyuan Wei, Jing Sun, Zijiang Zhang, Xianhao Zhang. **LLM-SmartAudit: Advanced Smart Contract Vulnerability Detection**. arXiv preprint arXiv:2410.09381 |
| [S27] | Mohamad Fakih, Rahul Dharmaji, Yasamin Moghaddas, Gustavo Quiros Araya, Oluwatosin Ogundare, Mohammad Abdullah Al Faruque. **LLM4PLC: Harnessing Large Language Models for Verifiable Programming of PLCs in Industrial Control Systems**. In Proceedings of the 46th International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP), 192 - 203 |
| [S28] | Wei Tao, Yucheng Zhou, Yanlin Wang, Wenqiang Zhang, Hongyu Zhang, Yu Cheng. **MAGIS: LLM-Based Multi-Agent Framework for GitHub Issue ReSolution**. arXiv preprint arXiv:2403.17927 |
| [S29] | Md. Ashraful Islam, Mohammed Eunus Ali, Md Rizwan Parvez. **MapCoder: Multi-Agent Code Generation for Competitive Problem Solving**. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL), 4912 - 4944 |
| [S30] | Dongming Jin, Zhi Jin, Xiaohong Chen, Chunhui Wang. **MARE: Multi-Agents Collaboration Framework for Requirements Engineering**. arXiv preprint arXiv:2405.03256 |
| [S31] | Yizhou Liu, Pengfei Gao, Xinchen Wang, Jie Liu, Yexuan Shi, Zhao Zhang, Chao Peng. **MarsCode Agent: AI-native Automated Bug Fixing**. arXiv preprint arXiv:2409.00899 |
| [S32] | Daman Arora, Atharv Sonwane, Nalin Wadhwa, Abhav Mehrotra, Saiteja Utpala, Ramakrishna Bairi, Aditya Kanade, Nagarajan Natarajan. **MASAI: Modular Architecture for Software-engineering AI Agents**. arXiv preprint arXiv:2406.11638 |
| [S33] | Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, Deheng Ye. **More Agents Is All You Need**. arXiv preprint arXiv:2402.05120 |
| [S34] | Zhuoyun Du, Chen Qian, Wei Liu, Zihao Xie, Yifei Wang, Yufan Dang, Weize Chen, Cheng Yang. **Multi-Agent Software Development through Cross-Team Collaboration**. arXiv preprint arXiv:2406.08979 |
| [S35] | Zhenyu Mao, Jialong Li, Dongming Jin, Munan Li, Kenji Tei. **Multi-role Consensus through LLMs Discussions for Vulnerability Detection**. arXiv preprint arXiv:2403.14274 |
| [S36] | Haolin Jin, Zechao Sun, Huaming Chen. **RGD: Multi-LLM Based Agent Debugger via Refinement and Generation Guidance**. arXiv preprint arXiv:2410.01242 |
| [S37] | Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, Zhiyuan Liu, Maosong Sun. **Scaling Large-Language-Model-based Multi-Agent Collaboration**. arXiv preprint arXiv:2406.07155 |
| [S38] | Yoichi Ishibashi, Yoshimasa Nishimura. **Self-Organized Agents: A LLM Multi-Agent Framework toward Ultra Large-Scale Code Generation and Optimization**. arXiv preprint arXiv:2404.02183 |
| [S39] | Haifeng Ruan, Yuntong Zhang, Abhik Roychoudhury. **SpecRover: Code Intent Extraction via LLMs**. arXiv preprint arXiv:2408.02232 |
| [S40] | John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, Ofir Press. **SWE-AGENT: AGENT-COMPUTER INTERFACES ENABLE AUTOMATED SOFTWARE ENGINEERING**. arXiv preprint arXiv:2405.15793 |
| [S41] | Noble Saji Mathews, Meiyappan Nagappan. **Test-Driven Development for Code Generation**. arXiv preprint arXiv:2402.13521 |
| [S42] | Feng Lin, Dong Jae Kim, Tse-Husn (Peter)Chen. **When LLM-based Code Generation Meets the Software Development Process**. arXiv preprint arXiv:2403.15852 |
| [S43] | Zhitao Wang, Wei Wang, Zirao Li, Long Wang, Can Yi, Xinjie Xu, Luyang Cao, Hanjing Su, Shouzhi Chen, Jun Zhou. **XUAT-Copilot: Multi-Agent Collaborative System for Automated User Acceptance Testing with Large Language Model**. arXiv preprint arXiv:2401.02705 |
| [S44] | Xunzhu Tang, Kisub Kim, Yewei Song, Cedric Lothritz, Bei Li, Saad Ezzini, Haoye Tian, Jacques Klein, Tegawende F. Bissyande. **CodeAgent: Autonomous Communicative Agents for Code Review**. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP), 11279 - 11313. |

| ID | Author, Publication Title, and Venue |
| --- | --- |
| [S45] | Dawen Zhang, Xiwei Xu, Chen Wang, Zhenchang Xing, Robert Mao. **A Layered Architecture for Developing and Enhancing Capabilities in Large Language Model-based Software Systems**. arXiv preprint arXiv:2411.12357 |
| [S46] | Sai Zhang, Zhenchang Xing, Ronghui Guo, Fangzhou Xu, Lei Chen, Zhaoyuan Zhang, Xiaowang Zhang, Zhiyong Feng, Zhiqiang Zhuang. **Empowering Agile-Based Generative Software Development through Human-AI Teamwork**. ACM Transactions on Software Engineering and Methodology, Volume 34, Issue 6, 1 - 46. |
| [S47] | Juyeon Yoon; Robert Feldt; Shin Yoo. **Intent-Driven Mobile GUI Testing with Autonomous Large Language Model Agents**. In Proceedings of the 17th IEEE International Conference on Software Testing, Verification & Validation (ICST), 129 - 139. |
| [S48] | Mohammadmehdi Ataei, Hyunmin Cheong, Daniele Grandi, Ye Wang, Nigel Morris, Alexander Tessier. **Elicitron: An LLM Agent-Based Simulation Framework for Design Requirements Elicitation**. arXiv preprint arXiv:2404.16045 |
| [S49] | Zhe Liu, Cheng Li, Chunyang Chen, Junjie Wang, Boyu Wu, Yawen Wang, Jun Hu, Qing Wang. **Vision-driven Automated Mobile GUI Testing via Multimodal Large Language Model**. arXiv preprint arXiv:2407.03037 |
| [S50] | Jiahong Xiang, Xiaoyang Xu, Fanchu Kong, Mingyuan Wu, Zizheng Zhang, Haotian Zhang, Yuqun Zhang. **How Far Can We Go with Practical Function-Level Program Repair?**. arXiv preprint arXiv:2404.12833 |
| [S51] | Zixiao Zhao, Jing Sun, Zhiyuan Wei, Cheng-Hao Cai, Zhe Hou, Jin Song Dong. **VisionCoder: Empowering Multi-Agent Auto-Programming for Image Processing with Hybrid LLMs**. arXiv preprint arXiv:2410.19245 |
| [S52] | Ahmed R. Sadik, Sebastian Brulin, Markus Olhofer, Antonello Ceravola, Frank Joublin. **LLM as a code generator in Agile Model Driven Development**. arXiv preprint arXiv:2410.18489 |
| [S53] | Yue Hu, Yuzhu Cai, Yaxin Du, Xinyu Zhu, Xiangrui Liu, Zijie Yu, Yuchen Hou, Shuo Tang, Siheng Chen. **Self-Evolving Multi-Agent Collaboration Networks for Software Development**. arXiv preprint arXiv:2410.16946 |
| [S54] | Zihan Liu, Ruinan Zeng, Dongxia Wang, Gengyun Peng, Jingyi Wang, Qiang Liu, Peiyu Liu, Wenhai Wang. **Agents4PLC: Automating Closed-loop PLC Code Generation and Verification in Industrial Control Systems using LLM-based Agents**. arXiv preprint arXiv:2410.14209 |
| [S55] | Xuanming Zhang, Yuxuan Chen, Yuan Yuan, Minlie Huang. **Seeker: Enhancing Exception Handling in Code with LLM-based Multi-Agent Approach**. arXiv preprint arXiv:2410.06949 |
| [S56] | Zhiqiang Yuan, Weitong Chen, Hanlin Wang, Kai Yu, Xin Peng, Yiling Lou. **TRANSAGENT: An LLM-Based Multi-Agent System for Code Translation**. arXiv preprint arXiv:2409.19894 |
| [S57] | Leilei Lin, Yingming Zhou, Wenlong Chen, Chen Qian. **Think-on-Process: Dynamic Process Generation for Collaborative Development of Multi-Agent System**. arXiv preprint arXiv:2409.06568 |
| [S58] | Huan Zhang, Wei Cheng, Yuhan Wu, Wei Hu. **A Pair Programming Framework for Code Generation via Multi-Plan Exploration and Feedback-Driven Refinement**. In Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering (ASE), 1319 - 1331. |
| [S59] | Weiqing Yang, Hanbin Wang, Zhenghao Liu, Xinze Li, Yukun Yan, Shuo Wang, Yu Gu, Minghe Yu, Zhiyuan Liu, Ge Yu. **Enhancing the Code Debugging Ability of LLMs via Communicative Agent Based Data Refinement**. arXiv preprint arXiv:2408.05006 |
| [S60] | Forough Mehralian, Titus Barik, Jeff Nichols, Amanda Swearngin. **Automated Code Fix Suggestions for Accessibility Issues in Mobile Apps**. arXiv preprint arXiv:2408.03827 |
| [S61] | Dong Huang, Jie M.Zhang, Michael Luck, Qingwen Bu, Yuhao Qing, Heming Cui. **AgentCoder: Multi-Agent Code Generation with Effective Testing and Self-optimisation**. arXiv preprint arXiv:2312.13010 |
| [S62] | Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, Jie Zhou. **AgentVerse: Facilitating Multi-Agent Collaboration and Exploring Emergent Behaviors**. In Proceedings of 17th International Conference on Learning Representations (ICLR). |
| [S63] | Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Jaward Sesay, Börje F. Karlsson, Jie Fu, Yemin Shi. **AutoAgents: A Framework for Automatic Agent Generation**. arXiv preprint arXiv:2309.17288 |
| [S64] | Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, Chi Wang. **AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation**. arXiv preprint arXiv:2308.08155 |
| [S65] | Zeeshan Rasheed, Muhammad Waseem, Kai-Kristian Kemell, Wang Xiaofeng, Anh Nguyen Duc, Kari Systä, Pekka Abrahamsson. **Autonomous Agents in Software Development: A Vision Paper**. arXiv preprint arXiv:2311.18440 |
| [S66] | Maryam Taeb, Amanda Swearngin, Eldon Schoop, Ruijia Cheng, Yue Jiang, Jeffrey Nichols. **AXNav: Replaying Accessibility Tests from Natural Language**. In Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems (CHI), 1 - 16. |
| [S67] | Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, Bernard Ghanem. **CAMEL: communicative agents for "mind" exploration of large language model society**. In Proceedings of the 37th International Conference on Neural Information Processing Systems (NeurIPS), 51991 - 52008 |
| [S68] | Guang Yang, Yu Zhou, Xiang Chen, Xiangyu Zhang, Terry Yue Zhuo, Taolue Chen. **Chain-of-Thought in Neural Code Generation: From and For Lightweight Language Models**. IEEE Transactions on Software Engineering, Volume 50, Issue 9, 2437 - 2457. |
| [S69] | Seungjun Moon, Hyungjoo Chae, Yongho Song, Taeyoon Kwon, Dongjin Kang, Kai Tzu-iunn Ong, Seung-won Hwang, Jinyoung Yeo. **Coffee: Boost Your Code LLMs by Fixing Bugs with Feedback**. arXiv preprint arXiv:2307.07924 |
| [S70] | Nalin Wadhwa, Jui Pradhan, Atharv Sonwane, Surya Prakash Sahu, Nagarajan Natarajan, Aditya Kanade, Suresh Parthasarathy, Sriram Rajamani. **CORE: Resolving Code Quality Issues using LLMs**. In Proceedings of the 31th ACM on Software Engineering (FSE), 789 - 811. |
| [S71] | Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, Diyi Yang. **A Dynamic LLM-Agent Network: An LLM-agent Collaboration Framework with Agent Team Optimization**. arXiv preprint arXiv:2310.02170. |
| [S72] | Yu Hao, Weiteng Chen, Ziqiao Zhou, Weidong Cui. **E&V: Prompting Large Language Models to Perform Static Analysis by Pseudo-code Execution and Verification**. arXiv preprint arXiv:2312.08477 |
| [S73] | Chen Qian, Yufan Dang, Jiahao Li, Wei Liu, Zihao Xie, Yifei Wang, Weize Chen, Cheng Yang, Xin Cong, Xiaoyin Che, Zhiyuan Liu, Maosong Sun. **Experiential Co-Learning of Software-Developing Agents**. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL), 5628 - 5640 |

| ID | Author, Publication Title, and Venue |
|---|---|
| [S74] | Martin Josifoski, Lars Klein, Maxime Peyrard, Nicolas Baldwin, Yifei Li, Saibo Geng, Julian Paul Schnitzler, Yuxing Yao, Jiheng Wei, Debjit Paul, Robert West. **Flows: Building Blocks of Reasoning and Collaborating AI**. arXiv preprint arXiv:2308.01285 |
| [S75] | Chunqiu Steven Xia, Matteo Paltenghi, Jia Le Tian, Michael Pradel, Lingming Zhang. **Fuzz4All: Universal Fuzzing with Large Language Models**. In Proceedings of the IEEE/ACM 46th International Conference on Software Engineering (ICSE), 1 - 13. |
| [S76] | Binfeng Xu, Xukun Liu, Hua Shen, Zeyu Han, Yuhan Li, Murong Yue, Zhiyuan Peng, Yuchen Liu, Ziyu Yao, Dongkuan Xu. **Gentopia.AI: A Collaborative Platform for Tool-Augmented LLMs**. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations (EMNLP), 237 - 245. |
| [S77] | Zhou Yang, Zhipeng Zhao, Chenyu Wang, Jieke Shi, Dongsum Kim, Donggyun Han, David Lo. **Gotcha! This Model Uses My Code! Evaluating Membership Leakage Risks in Code Models**. arXiv preprint arXiv:2310.01166 |
| [S78] | Hanbin Wang, Zhenghao Liu, Shuo Wang, Ganqu Cui, Ning Ding, Zhiyuan Liu, Ge Yu. **INTERVENOR: Prompting the Coding Ability of Large Language Models with the Interactive Chain of Repair**. In Proceedings of the 62nd Findings of the Association for Computational Linguistics (ACL), 2081 - 2107. |
| [S79] | 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, Tao Yu. **Lemur: Harmonizing Natural Language and Code for Language Agents**. In Proceedings of the 12th International Conference on Representation Learning (ICLR). |
| [S80] | Zhe Liu, Chunyang Chen, Junjie Wang, Mengzhuo Chen, Boyu Wu, Xing Che, Dandan Wang, Qing Wang. **Make LLM a Testing Expert: Bringing Human-like Interaction to Mobile GUI Testing via Functionality-aware Decisions**. In Proceedings of the IEEE/ACM 46th International Conference on Software Engineering (ICSE), 1 - 13. |
| [S81] | 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, Jürgen Schmidhuber. **MetaGPT: Meta Programming for A Multi-Agent Collaborative Framework**. arXiv preprint arXiv:2308.00352. |
| [S82] | Yashar Talebirad, Amirhossein Nadiri. **Multi-Agent Collaboration: Harnessing the Power of Intelligent LLM Agents**. arXiv preprint arXiv:2306.03314. |
| [S83] | Eric Zelikman, Qian Huang, Gabriel Poesia, Noah D. Goodman, Nick Haber. **Parsel: Algorithmic Reasoning with Language Models by Composing Decompositions**. In Proceedings of 37th Conference on Neural Information Processing Systems (NeurIPS). |
| [S84] | Zhenchang Xing, Qing Huang, Yu Cheng, Liming Zhu, Qinghua Lu, Xiwei Xu. **Prompt Sapper: LLM-Empowered Software Engineering Infrastructure for AI-Native Services**. arXiv preprint arXiv:2306.02230 |
| [S85] | Zefan Wang, Zichuan Liu, Yingying Zhang, Aoxiao Zhong, Jihong Wang, Fengbin Yin, Lunting Fan, Lingfei Wu, Qingsong Wen. **RCAgent: Cloud Root Cause Analysis by Autonomous Agents with Tool-Augmented Large Language Models**. arXiv preprint arXiv:2310.16340 |
| [S86] | Yihong Dong, Xue Jiang, Zhi Jin, Ge Li. **Self-collaboration Code Generation via ChatGPT**. arXiv preprint arXiv:2304.07590 |
| [S87] | Kechi Zhang, Zhuo Li, Jia Li, Ge Li, Zhi Jin. **Self-Edit: Fault-Aware Code Editor for Code Generation**. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL), 769 - 787. |
| [S88] | Gang Fan, Xiaoheng Xie, Xunjin Zheng, Yinan Liang, Peng Di. **Static Code Analysis in the AI Era: An In-depth Exploration of the Concept, Function, and Potential of Intelligent Code Analysis**. arXiv preprint arXiv:2310.08837 |
| [S89] | Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, Karthik Narasimhan. **Swe-bench: Can language models resolve real-world github issues?**. In Proceedings of the 12th International Conference on Representation Learning (ICLR). |
| [S90] | Xinyun Chen, Maxwell Lin, Nathanael Schärli, Denny Zhou. **Teaching Large Language Models to Self-Debug**. In Proceedings of the 12th International Conference on Representation Learning (ICLR). |
| [S91] | Chenyuan Yang, Yinlin Deng, Runyu Lu, Jiayi Yao, Jiawei Liu, Reyhaneh Jabbarvand, Lingming Zhang. **White-box Compiler FuzzingEmpowered by Large Language Models**. arXiv preprint arXiv:2310.15991. |
| [S92] | Theo X. Olausson, Jeevana Priya Inala, Chenglong Wang, Jianfeng Gao, Armando Solar-Lezama. **Is Self-Repair a Silver Bullet for Code Generation?**. In Proceedings of the 12th International Conference on Representation Learning (ICLR). |
| [S93] | Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, Maosong Sun. **ChatDev: Communicative Agents for Software Development**. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL), 15174 - 15186. |
| [S94] | Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Vincent Zhao, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Pranesh Srinivasan, Laichee Man, Kathleen Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguera-Arcas, Claire Cui, Marian Croak, Ed Chi, Quoc Le. **LaMDA: Language Models for Dialog Applications**. arXiv preprint arXiv:2201.08239. |