

# A Survey on the Optimization of Large Language Model-based Agents

SHANGHENG DU, Shanghai Institute of Artificial Intelligence for Education, East China Normal University; School of Computer Science and Technology, East China Normal University, China  
 JIABAO ZHAO\*, School of Computer Science and Technology, Donghua University, China  
 JINXIN SHI, Shanghai Institute of Artificial Intelligence for Education, East China Normal University; School of Computer Science and Technology, East China Normal University, China  
 ZHENTAO XIE, School of Computer Science and Technology, East China Normal University, China  
 XIN JIANG, School of Computer Science and Technology, East China Normal University, China  
 YANHONG BAI, Shanghai Institute of Artificial Intelligence for Education, East China Normal University; School of Computer Science and Technology, East China Normal University, China  
 LIANG HE, School of Computer Science and Technology, East China Normal University, China

With the rapid development of Large Language Models (LLMs), LLM-based agents have been widely adopted in various fields, becoming essential for autonomous decision-making and interactive tasks. However, current work typically relies on prompt design or fine-tuning strategies applied to vanilla LLMs, which often leads to limited effectiveness or suboptimal performance in complex agent-related environments. Although LLM optimization techniques can improve model performance across many general tasks, they lack specialized optimization towards critical agent functionalities such as long-term planning, dynamic environmental interaction, and complex decision-making. Although numerous recent studies have explored various strategies to optimize LLM-based agents for complex agent tasks, a systematic review summarizing and comparing these methods from a holistic perspective is still lacking. In this survey, we provide a comprehensive review of LLM-based agent optimization approaches, categorizing them into parameter-driven and parameter-free methods. We first focus on parameter-driven optimization, covering fine-tuning-based optimization, reinforcement learning-based optimization, and hybrid strategies, analyzing key aspects such as trajectory data construction, fine-tuning techniques, reward function design, and optimization algorithms. Additionally, we briefly discuss parameter-free strategies that optimize agent behavior through prompt engineering and external knowledge retrieval. Finally, we summarize the datasets and benchmarks used for evaluation and tuning, review key applications of LLM-based agents, and discuss major challenges and promising future directions. Our repository for related references is available at <https://github.com/YoungDubbyDu/LLM-Agent-Optimization>.

## 1 Introduction

The development of autonomous agents has been a long-term pursuit in Artificial Intelligence (AI). AI agents have evolved from early rule-based and expert system-based architectures to reinforcement learning (RL)-driven agents, which are now widely applied in many fields [35]. Traditional RL-based agents optimize policies through interaction with environments, using structured reward functions to achieve goals and improve performance over time. However, these approaches often

\*Corresponding author.

Authors' Contact Information: Shangheng Du, Shanghai Institute of Artificial Intelligence for Education, East China Normal University; School of Computer Science and Technology, East China Normal University, Shanghai, China, dsh@stu.ecnu.edu.cn; Jiabao Zhao, School of Computer Science and Technology, Donghua University, Shanghai, China, jbzhaod@dhru.edu.cn; Jinxin Shi, Shanghai Institute of Artificial Intelligence for Education, East China Normal University; School of Computer Science and Technology, East China Normal University, Shanghai, China, jinxinshi@stu.ecnu.edu.cn; Zhentao Xie, School of Computer Science and Technology, East China Normal University, Shanghai, China, ecnudavidtao@gmail.com; Xin Jiang, School of Computer Science and Technology, East China Normal University, Shanghai, China, 51275901099@stu.ecnu.edu.cn; Yanhong Bai, Shanghai Institute of Artificial Intelligence for Education, East China Normal University; School of Computer Science and Technology, East China Normal University, Shanghai, China, Lucky\_Baiyh@stu.ecnu.edu.cn; Liang He, School of Computer Science and Technology, East China Normal University, Shanghai, China, lhe@cs.ecnu.edu.cn.

require extensive training, rely on well-defined state-action spaces, and struggle with generalization across diverse tasks.

In recent years, Large Language Models (LLMs) such as GPT-4 [120], PaLM 2 [5], and Deepseek-r1 [52] have achieved remarkable success, demonstrating exceptional capabilities in language understanding, reasoning, planning and complex decision-making. Building on these strengths, LLMs can serve as agents, providing a promising pathway to improve autonomous decision-making and achieve AGI [169]. Unlike conventional RL-based agents, which optimize explicit reward-driven policies, LLM-based agents operate through text-based instructions and prompt templates and in-context learning (ICL), allowing greater flexibility and generalization. These agents leverage the comprehension and reasoning capabilities of LLMs to interact with environments through natural language, execute complex multi-step tasks, and dynamically adapt to evolving scenarios. Existing LLM agents utilize various methods such as task decomposition [64], self-reflection [133], memory augmentation [210], and multi-agent collaboration [86] to achieve high performance across a range of domains, including software development [67], mathematical reasoning [1], embodied intelligence [212], web navigation [28], and more.

However, despite their strengths, LLMs are not inherently designed for autonomous decision-making and long-term tasks. Their training objectives focus on next-token prediction rather than reasoning, planning, or interactive learning required for agent-based tasks, so they lack explicit training on agent-centric tasks. As a result, deploying LLMs as agents in complex environments presents several key challenges: 1) LLM-based agents struggle with long-horizon planning and multi-step reasoning, as their generative content may lead to task inconsistencies or error accumulation over extended interactions. 2) Limited memory capacity in LLMs hinders agents from utilizing past experiences for reflection, leading to suboptimal decision-making and task performance. 3) The adaptability of LLM-based agents to novel environments is constrained, as they primarily rely on pre-trained knowledge or fixed contexts, limiting their ability to handle dynamic scenarios. These limitations are particularly evident in open-source LLMs, which lag behind proprietary models like GPT-4 in agent-specific capabilities. Additionally, the high cost and lack of transparency of closed-source LLMs highlight the need for optimizing open LLMs to enhance agent capabilities.

Existing techniques, such as supervised fine-tuning (SFT) [122] and reinforcement learning with human feedback (RLHF) [121], have made significant strides in improving LLM performance in instruction following tasks, but they fail to fully address the challenges of decision-making, long-term planning, and adaptability for LLM-based agents. Optimizing LLM-based agents requires a broader understanding of dynamic environments and agent behaviors, which needs to design specialized techniques that go beyond traditional LLM fine-tuning and prompt engineering methods. To address these challenges, numerous recent studies have explored various strategies to optimize LLM-based agents for complex agent tasks. These methods ensure that agents can generalize across diverse environments, refine strategies based on feedback, and efficiently utilize external resources such as tools, memory, and retrieval mechanisms.

In this paper, we provide a comprehensive survey on LLM-based agent optimization, systematically categorizing methods into parameter-driven and parameter-free optimization strategies. Our work focuses on the technical methodologies employed to optimize agent capabilities like agent tuning, RL, and others to improve agent performance. Specifically, **Parameter-driven Optimization** refines LLM parameters to enhance agent performance. This category includes conventional fine-tuning approaches, covering key stages such as agent trajectory data construction and fine-tuning strategies. In addition, we explore RL-based optimization, which is divided into two distinct optimization directions: reward function-based methods leveraging traditional RL techniques like Actor-Critic [147] and Proximal Policy Optimization (PPO) [136], and preference alignment-based methods utilizing Direct Preference Optimization (DPO) [132] to synchronize

agent policies with human preference or task-specific objectives. Finally, we discuss hybrid fine-tuning optimization strategies, a rising area, which combine SFT with RL to iteratively refine agent behavior. In contrast, we also briefly outline **Parameter-free Optimization** methods that focus on improving agent behavior without modifying model parameters. These methods leverage prompt engineering, in-context learning and retrieval-augmented generation (RAG), incorporating various types of information into prompts to guide agents' actions. They are categorized into feedback-based optimization, experience-based optimization, tool-based optimization, retrieval-augmented optimization, and multi-agent collaborative optimization.

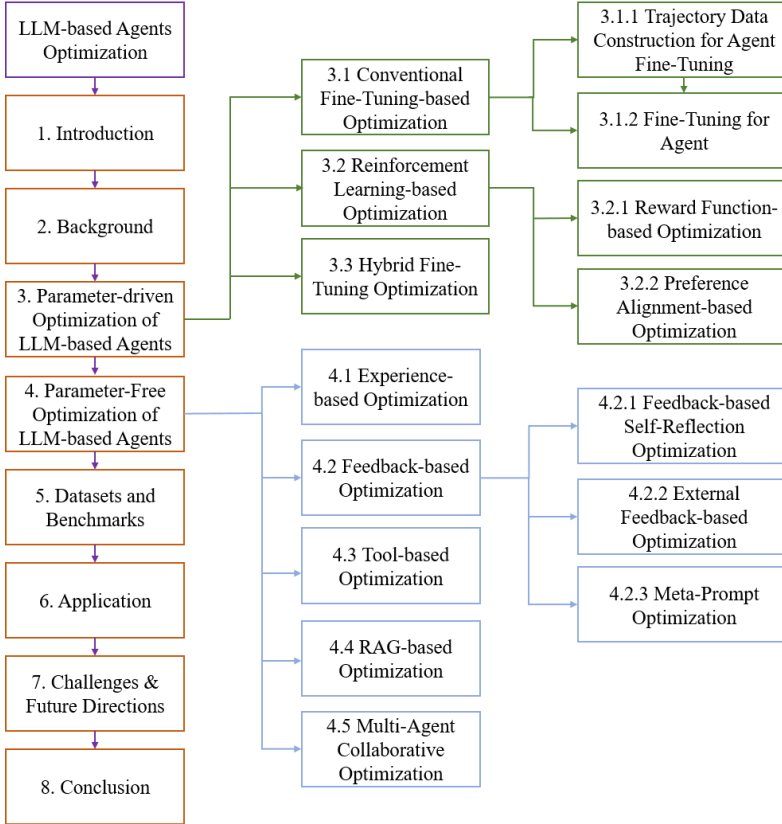


Fig. 1. An Overview of the Paper Organization.

**Comparison to related surveys.** Despite the growing research interest in LLM-based agents, existing surveys primarily focus on general LLM optimization or specific agent abilities such as planning, memory, and role-playing, without treating LLM-based agent optimization as a distinct research area. Surveys on LLM optimization mainly cover fine-tuning [115, 122] and self-evolution approaches [150], but lack discussions on specialized optimization required for agent capabilities. On the other hand, existing agent-related surveys generally categorize works based on architectural components such as planning [64], memory [210], or multi-agent coordination [86], rather than systematically summarizing the techniques dedicated to optimize LLM-based agent behaviors and performance. As comparison, this work is the first survey towards LLM-based agent optimization

techniques, facilitating a clearer understanding and comparison of existing methods and providing directions for future research.

**Scope and rationales.** (1) We survey only LLM-based agent optimization algorithms to improve agent task performance, such as problem-solving and decision-making, covering parameter-driven and parameter-free approaches. We exclude works centered on general LLM efficiency, role-playing, or dialogue; (2) Our selection includes papers from AI and NLP conferences and journals, as well as recent high-impact preprints from arXiv to ensure coverage of the latest advancements. (3) We focus on studies published since 2022 to reflect recent advancements in LLM-based agent optimization.

**Organization of this survey.** The schematic representation of this manuscript’s layout can be found in Figure 1. Section 2 provides the background knowledge and related concepts. In Section 3, we systematically review parameter-driven optimization approaches that modify LLM parameters to enhance agent capabilities, categorizing them into three main strategies: fine-tuning-based optimization (§3.1), RL-based optimization (§3.2), and hybrid optimization (§3.3). Section 4 summarizes and classifies existing work on parameter-free optimization strategies. Then, Section 5 presents datasets and benchmarks, while Section 6 reviews practical applications across various domains. Finally, Section 7 highlights challenges and future directions.

## 2 Background

### 2.1 Reinforcement Learning-based Agent Optimization

RL has long been a fundamental approach in agent optimization, allowing agents to learn from interactions with environments. Current RL methods mainly optimize agent behaviors using value-based and policy-based approaches [35, 106, 117]. **Value-based methods**, such as Q-learning [25, 163], optimize an agent’s action-value function to maximize long-term rewards. These methods are effective in discrete action spaces but struggle with high-dimensional states or action spaces. **Policy-based methods**, including Policy Gradient [48, 124], directly optimize the agent’s policy by adjusting parameters based on reward gradients. To improve stability and sample efficiency, PPO [136] introduced a constraint on policy updates, mitigating performance degradation during training. Actor-Critic methods [147] combine value estimation with policy learning, improving convergence efficiency and decision robustness. Beyond single-agent settings, Multi-Agent Reinforcement Learning (MARL) extends RL techniques to scenarios involving multiple interacting agents, enabling both cooperative and competitive dynamics [12, 204].

In recent years, RL has also been increasingly applied to aligning AI agents with human intentions, particularly in preference-based optimization. RLHF [121] has emerged as a prominent approach, refining agent policies based on human-provided signals to improve alignment with desired behaviors. DPO [132] optimizes policies directly from preference data without reward modeling, improving alignment and controllability. Overall, RL-based optimization has evolved from early value-based and policy-based learning to more advanced techniques that integrate structured feedback and multi-agent coordination, providing a foundation for improving decision-making in LLM-based agents.

### 2.2 LLM Fine-Tuning

LLM fine-tuning is a critical method for adapting pre-trained models to specific tasks through optimizing parameters, making them more suited to the desired application. The most popular approach is SFT, where LLMs are trained on labeled data to improve task-specific performance. Instruction Tuning is a commonly used method in SFT, where LLMs are further trained on instruction-output pairs to enhance their ability to follow human commands [98, 205]. Another major development is parameter-efficient fine-tuning (PEFT), including methods like P-Tuning [103], LoRA [59], and

QLoRA [30]. These techniques adjust a small subset of parameters, significantly reducing the computational cost of fine-tuning while preserving LLM performance, making them highly efficient for real-world applications. Additionally, RLHF has been used to fine-tune LLMs by integrating human feedback, improving their decision-making and output alignment with user preferences [121]. These optimization techniques enable LLMs to adapt more efficiently to a wide range of tasks, enhancing their effectiveness in real-world scenarios.

### 2.3 LLM-based RAG

RAG combines LLM with external information retrieval systems to enhance the relevance and accuracy of generated outputs. By retrieving relevant documents from external sources, RAG allows LLMs to address the knowledge constraints inherent in models. The evolution of RAG methods has been marked by significant advancements in retrieval and generation integration [44]. Early, Naive RAG methods focus on directly retrieving relevant documents to augment the generative process, improving the quality of responses in tasks requiring factual knowledge. To address the challenges of Naive RAG, Advanced RAG is introduced, refining the retrieval process by incorporating more effective ranking, filtering, and document selection strategies. Subsequently, Modular RAG introduces a modular framework that optimizes the retrieval and generative components independently. This modular approach enables task-specific optimizations, allowing for more flexibility and scalability in applications across different domains [8, 193]. These advancements in RAG highlight its potential to enhance LLMs by enabling dynamic access to external knowledge, making them more adaptable and capable of addressing complex tasks in real-world scenarios.

## 3 Parameter-driven Optimization of LLM-based Agents

**Comparison with LLM parameter optimization.** Parameter-driven LLM optimization focuses on "how to create a better model", aiming to enhance general language understanding, instruction following, and broad task performance. In contrast, LLM-based agent parameter optimization addresses "how to use the model to solve complex agent tasks", emphasizing decision-making, multi-step reasoning, and task execution in dynamic environments. Although general LLM optimization improves fluency and factual accuracy across diverse applications, LLM-agent optimization is task-specific, requiring models to adapt strategies, interact with environments, and refine behaviors for autonomous problem-solving. Parameter-driven optimization of LLM-based agents primarily relies on expert trajectory data or self-generated trajectory data obtained through environment exploration, then employs various optimization techniques to iteratively refine policies and enhance performance.

In this section, we discuss how parameter-driven optimization methods improve the performance of LLM-based agents. Specifically, we categorize these methods into three main technical approaches according to different strategies for parameter tuning: conventional fine-tuning-based optimization, reinforcement learning-based optimization, and hybrid optimization.

### 3.1 Conventional Fine-Tuning-based Optimization

Conventional fine-tuning-based agent optimization involves tuning pre-trained LLMs' parameters through various fine-tuning techniques, such as instruction tuning and parameter-efficient fine-tuning. Trajectory for fine-tuning typically are constructed in the form of SFT and is used to adjust the agent's parameters to better align with task-specific requirements. The optimization process typically consists of two major steps: **1) constructing high-quality trajectory data tailored to agent tasks; 2) fine-tuning LLM-based agents using these trajectory data**, and the complete process is presented in Figure 2. Previous studies [40, 83, 122] have shown that the quality of training data significantly impacts model performance, highlighting the importance

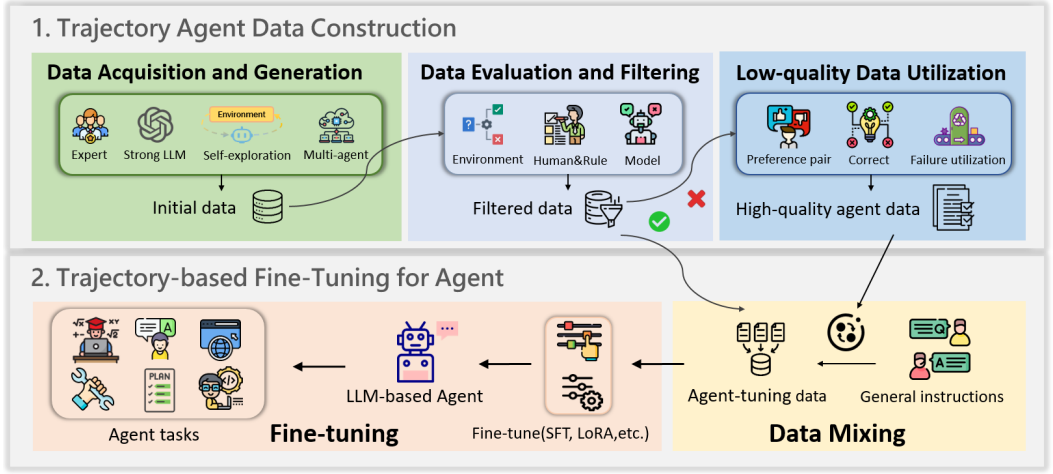


Fig. 2. Workflow of Fine-Tuning-based Optimization for LLM-based Agents.

of generating, filtering, and effectively utilizing high-quality trajectories. This makes trajectory construction a critical step in the fine-tuning pipeline, directly influencing the LLM-based agent’s overall performance. In Table 1, we provide a comprehensive overview of fine-tuning-based agent optimization methods, highlighting the data processing techniques and fine-tuning strategies used in each work. It is important to note that this section excludes fine-tuning methods that involve reinforcement learning or preference alignment techniques (e.g., DPO, PPO), which will be addressed in §3.2. Instead, in this section, we only focus on the part of traditional LLM fine-tuning techniques applied in existing works, aiming to ensure each stage of the conventional fine-tuning-based agent optimization workflow is clearly introduced.

**3.1.1 Trajectory Data Construction for Agent Fine-Tuning.** The construction of high-quality trajectory is a crucial step before the fine-tuning of LLM-based agents, which aims to empower LLMs with agent ability. This process involves the generation of trajectory data, followed by evaluation and filtering, and the potential utilization of low-quality samples, to construct refined data that meet the requirements for effective fine-tuning.

**Data Acquisition and Generation.** High-quality trajectory data construction begins with the acquisition and generation of initial data, which requires not only a diverse set of trajectories, but also sufficient alignment with the target tasks to ensure effective learning. Methods for acquiring and generating such data can generally be classified into four broad categories: expert-annotated data, strong LLM-generated trajectories, self-exploration environment-interaction trajectories, and multi-agent collaboration-based construction. Here, we introduce the utilization and construction processes of each category and review the relevant studies.

**(1) Expert-annotated data.** Expert-annotated trajectories refer to high-quality datasets manually crafted by human experts, often considered the gold standard for fine-tuning. These data ensure task reliability and alignment, as experts can meticulously design and annotate trajectories tailored to specific cases.

Many works [14, 39, 144, 158, 177] utilize ReAct-style expert trajectories as initial datasets, with data including thoughts, observations and actions [189], which enable agents to mimic expert decision-making processes more effectively. For instance, IPR [177] leverages such trajectories to help agents acquire foundational skills. Similarly, ETO [144] and AGILE [39] apply Chain of Thought

Table 1. Comparison of Conventional Fine-Tuning-based Optimization for LLM-based Agents: Data Construction and Fine-Tuning. Note: MA - Multi-Agent Framework; LQ - Low-Quality Data Utilization.

Method	Trajectory Agent Data Construction				Fine-Tuning	
	Generation	Filtering	MA	LQ	Fine-tune Approach	Base Model
AgentTuning [199]	Strong LLM	Human or Rule	✓	✓	Instruction Tuning	Llama-2-7B/13B/70B
SMART [197]	Multi-agent	Environment	/	✓	LoRA	Llama-2-7B
Agent-FLAN [22]	Expert	Model	✓	✓	Instruction Tuning	Llama-2-7B
Self-Talk [153]	Multi-agent	Human or Rule	/	✓	LoRA	MosaicAI-7B-Chat
ENVISIONS [178]	Self-exploration	Environment	/	✓	SFT	Llama2-7B/13B-Chat
AgentGym [170]	Strong LLM&Expert	Environment	/	✓	BC	Llama-2-7B-Chat
FireAct [14]	Strong LLM	Environment	/	/	LoRA	GPT3.5, Llama-2-7B/13B, CodeLlama-7B/13B/34B-Instruct
NAT [158]	Strong LLM	Environment	/	✓	SFT	Llama-2-7B/13B-Chat
Agent Lumos [192]	Strong LLM	Human or Rule	/	/	LoRA	Llama-2-7B/13B
STE [154]	Self-exploration	Model	/	✓	SFT	Llama-2-7B/13B-Chat, Mistral-7B-Instruct
OPTIMA [19]	Multi-agent	Human or Rule	✓	/	SFT	Llama-3-8B
Zhou et al. [216]	Strong LLM	Human or Rule	✓	/	LoRA	OpenChat v3.2, Llama-2-7B, AgentLM-7B
AgentOhana [202]	Expert	Model	/	/	QLoRA	xLAM-v0.1
COEVOL [85]	Expert	Model	✓	/	SFT	Llama-2-7B, Mistral-7B
AGENTBANK [143]	Strong LLM	Environment	/	✓	Instruction Tuning	Llama-2-Chat
ADASWITCH [146]	Self-exploration	Model	✓	✓	SFT	DeepSeek-Coder-1.3B, StarCoder2-3B
IPR [177]	Expert & Self-exploration	Environment	/	✓	Instruction Tuning	Llama-2-7B
Re-ReST [33]	Self-exploration	Environment	/	✓	LoRA	Llama-2-7B/13B, Llama-3-8B, CodeLlama-13B, VPGen
ATM [219]	Multi-agent	/	✓	/	MITO	Llama-2-7B
Aksitov et al. [3]	Self-exploration	Model-based	/	/	SFT	PaLM-2-base-series
SWIFTSAGE [94]	Self-exploration	Environment	✓	/	SFT	T5-Large
AGILE [39]	Expert	/	/	/	BC	Vicuna-13B, Meerkat-7B
NLRL [40]	Self-exploration	/	/	/	SFT	Llama-3.1-8B-Instruct
ETO [144]	Expert	/	/	✓	BC	Llama-2-7B-Chat
Retrospec [171]	Expert	/	/	✓	BC	Flan-T5-Large, Llama-3-8B-Instruct
ToRA [49]	Strong LLM	Human or Rule	/	✓	BC	Llama-2-series, CodeLlama-series
Sayself [179]	Strong LLM	Human or Rule	/	/	SFT	Mistral-7B, Llama-3-8B

(CoT) methods [164] to expert trajectories for imitation learning, reinforcing task-specific behaviors. To ensure alignment with pre-trained LLM domains, Agent-FLAN [22] transforms ReAct-style expert trajectories into multi-turn dialogue, segmenting the dialogue into different task-specific turn, such as instruction-following and reasoning. StepAgent [29] introduces a two-phase learning process, where agents first observe discrepancies between their policies and expert trajectories, then iteratively refine their actions. Additionally, AgentOhana [202] standardizes heterogeneous agent expert trajectories into a unified format to improve data consistency. Despite their reliability and alignment with specific tasks, these datasets are resource-intensive and lack scalability, making them commonly supplemented with other data acquisition methods to enhance dataset diversity.

**(2) Strong LLM-generated trajectories.** Strong LLM-generated trajectories leverage powerful LLMs like ChatGPT and GPT-4 to autonomously generate task-specific data. These trajectories are usually produced by reasoning frameworks such as ReAct and CoT, allowing the model to interact with the environment and simulate processes of reasoning, decision-making and acting.

AgentTuning [199] and FireAct [14] employ ReAct and CoT to guide agent behavior while incorporating Reflexion [139] refinements, improving the diversity of generated data. Some works integrate tools and structured annotations to enhance trajectory informativeness. NAT [158] generates multiple trajectories under different temperature settings, using ReAct prompts and integrating tools such as calculators and APIs during interactions. Agent Lumos [192] utilizes GPT-4 and GPT-4V to annotate datasets within planning and grounding modules, producing LUMOS-I and LUMOS-O style data. Other methods explore multi-role simulation to enrich trajectory complexity. Zhou et al. [216] employ GPT-4 to simulate problem generators, action planners, and environment agents, enabling iterative interaction-driven data generation. AGENTBANK [143] also leverages GPT-4 for environment interaction data and GPT-3.5 for CoT rationales, and finally transforms the data into chatbot-style formats for improved usability.

**(3) Self-exploration environment-interaction trajectories.** Given the high costs of expert annotation and proprietary models like GPT-4, self-exploration methods using open-source models

have become a common approach for generating trajectory data. These methods enable agents to interact with environments and generate trajectories through autonomous exploration, and feedback-based learning and self-training, thereby reducing the reliance on manual annotation.

Early approaches [3, 94] focus on direct exploration and interaction with the environment. For example, SWIFTSAGE [94] utilizes longer trajectory history representations to selectively sample actions and construct concise training datasets. Similarly, Aksitov et al. [3] applies ReAct agent, proposed by Reflexion [139], to optimize trajectories, ensuring better task alignment. Building on this foundation, recent methods incorporate feedback mechanisms and self-reflection to refine trajectory quality. ENVISIONS [178] allows agents to explore and refine their trajectories through self-correction and self-reward mechanisms. STE [154] simulates tool usage and API interactions, storing feedback and experiences as memory to enhance trajectory learning. Other approaches focus on leveraging additional models for reflection and correction. ADASWITCH [146] employs a hybrid setup where a local-deployed LLM executes tasks, and a cloud-based LLM corrects errors, integrating refined actions into structured datasets. Re-ReST [33] iteratively refines trajectory outputs using a Reflector model, aligning them more closely with task requirements. Additionally, NLRL [40] introduces an innovative perspective by using natural language to represent RL optimization processes, such as text-style policies, value functions, and feedback, enabling agents to interact with environments and generate trajectory data through RL view.

Table 2. Comparison of Data Acquisition and Generation Strategies.

Strategy	Advantages	Limitations
Expert-annotated	Ensures high accuracy with flexible, task-specific manual design.	Labor-intensive, limited, and difficult to obtain in professional domains.
Strong LLM-generated	Provides stable and high-performance trajectories without complex design.	Incurs high usage cost and may introduce potential biases and factual inconsistencies.
Self-exploration	Offers low cost, easy implementation, and no need for extra assistance.	Easily produces low-quality trajectories with errors, requiring effective data filtering.
Multi-agent collaboration	Improves accuracy and scalability through task specialization.	Relies on complex coordination mechanisms, which increase design complexity and cost.

**(4) Multi-agent collaboration-based construction.** Multi-agent collaborative frameworks overcome the limitations of single-agent methods by enabling richer interactions and coordinated dynamics, resulting in more diverse and robust trajectory data. Although these methods also often rely on strong LLMs or self-exploration with environments, the multi-agent design introduces unique mechanisms that emphasize collaborative dynamics and modular task execution, making them worthy of separate discussion.

Most multi-agent collaboration-based methods [85, 153, 197, 219] employ different roles for each agent, allowing them to collaboratively generate or refine trajectory data through sequential or alternating iterative processes. In SMART [197], agents sequentially play roles as intent reconstructor, knowledge retriever, fact locator, and response generator, collaboratively producing comprehensive answers and constructing trajectory datasets that encompass both short-term and long-term interactions. Similarly, COEVOL [85] employs a two-stage debate mechanism to generate initial data, and evaluate and refine them by agents acting as an advisor, editor, and judge to ensure high-quality outputs. Self-Talk [153] enables LLMs to simulate conversations among diverse roles, generating training data that reflects multiple perspectives. To enhance the robustness of data, ATM [219] incorporates adversarial interactions, with an attacker perturbing RAG-retrieved information while a generator produces responses that are robust to these disturbances. By dividing tasks among agents and facilitating sequential or cross refinement, multi-agent collaborative frameworks enhance trajectory generation by improving data granularity, diversity, and adaptability.



**Summary.** The data acquisition and generation strategies play a crucial role in optimizing LLM-based agents by offering different approaches to constructing trajectory data, each with unique advantages and limitations, as summarized in Table 2. Expert-annotated data ensure high reliability and accuracy but are labor-intensive and costly, making large-scale collection challenging. Consequently, they are often integrated with other data sources to balance quality and efficiency. Strong LLM-generated trajectories address scalability issues by leveraging powerful models like GPT-4, but their reliance on these advanced models leads to high costs in API usage or computational resources. Additionally, their quality depends on the intrinsic reasoning capabilities of the LLMs, requiring further refinement to mitigate model biases. Self-exploration methods enable agents to interact with environments and autonomously generate task-specific data, reducing reliance on human supervision or closed-source LLMs. However, due to the limited capabilities of models in exploration, they often produce a large number of low-quality or failed trajectories, requiring effective filtering to extract useful data for fine-tuning. Multi-agent collaboration frameworks enhance data diversity and adaptability through task distribution and iterative refinement. However, they rely heavily on effective coordination among agents, making the system structurally unstable and increasing design complexity. Additionally, the overall cost rises as more agents are involved. The suitability of each data acquisition method depends on specific task requirements, as different approaches offer varying trade-offs in quality, scalability, and complexity. Selecting the appropriate strategy should be guided by the research needs, and in the future, exploring hybrid strategies that integrate multiple methods could help leverage their respective strengths.

**Data Evaluation and Filtering.** After trajectory data is generated, evaluating and filtering the data becomes essential to ensure their quality and suitability for fine-tuning. Although expert-annotated datasets are often directly usable because of their high initial quality, other types of trajectories typically require further refinement. These processes involve evaluating trajectories based on predefined criteria, external feedback or model-based scoring, followed by filtering or correcting the data to align with specific fine-tuning objectives. We categorize data evaluation and filtering methods into three primary approaches: Environment-based, Human or Rule-based, and Model-based evaluation.

**(1) Environment-based.** Environment-based methods rely on external feedback from the environment to assess the quality and success of generated trajectories. These approaches typically utilize environmental rewards or task completion indicators to assess agents' actions, and filter and refine trajectory data. Most works [14, 170, 178, 199] in this category adopt a form of binary feedback, where the reward is assigned based on the success or failure of the agent's actions relative to the desired task. AgentTuning [199], ENVISIONS [178] and FireAct [14] use environment-based rewards to evaluate trajectory success, relying on environmental signals or correctness assessments. Similarly, NAT [158] also employs a binary environment-based reward system, differentiating between successful and unsuccessful outcomes to encourage agents to learn from failure.

**(2) Human or rule-based.** Human or rule-based evaluation methods rely on predefined criteria or custom rules to assess the quality of trajectory data. The rules can be based on various metrics, such as task performance, data diversity, or consistency, and can be adapted to meet the specific requirements of different tasks.

Many works [3, 192, 216] rely on human annotation or manual review to evaluate and filter trajectory data, ensuring accuracy and alignment with predefined standards. In [216], manual filtering is applied to ensure logical consistency between the environment agent's feedback and the actor's actions. Similarly, Agent Lumos [192] leverages ground truth for evaluation and filtering, and ensures that final responses meet the established standards for task performance. Perplexity (PPL) is employed as a filtering criterion, selecting the samples with the lowest perplexity to aligns with the task goals [3]. In contrast, some methods [19, 153, 178] introduce custom evaluation

criteria that consider multiple dimensions of the data. Self-Talk [153] uses automated metrics to evaluate data based on dialogue diversity, sub-goal completion, and role consistency through different filters. ENVISIONS [178] designs a reward function, ranking generated trajectories based on sequence output probabilities. OPTIMA [19] balances task performance, token efficiency, and natural language readability in its customized reward function, and selects the top-K trajectories.

**(3) Model-based.** Model-based approaches leverage pre-trained LLMs (e.g., GPT-3.5/4) to automate trajectory quality assessment through model-driven evaluation mechanisms. These methods employ models’ reasoning capabilities to perform multidimensional analyses of relevance, accuracy, and completeness. This ensures alignment of trajectories and task while establishing a holistic quality assurance framework for data construction.

Some approaches utilize the reasoning and scoring capabilities of LLMs to select the most relevant and correct trajectories directly. [3] introduces a ranking approach, using a fine-tuned PaLM 2-L model to sort multiple samples based on their results, selecting the highest-ranked samples for further processing. STE [154] uses GPT-4 to evaluate the effectiveness of using examples about APIs, simulated tool interactions. ADASWITCH [146] employs the LLM voting verifier or roberta-large-mnli [105] to check the correctness of each step and adjust accordingly. Other approaches focus on evaluating the entire trajectory rather than just the input-output pair. AgentOhana [202] employs the AgentRater system to evaluate trajectories on a 0-5 scale, assessing both process integrity and outcome validity, retaining only trajectories scoring  $\geq 4$  to ensure comprehensive logical consistency. Additionally, COEVOL [85] implements GPT-3.5 for multidimensional quality assessment, conducting comparative evaluations between raw and processed trajectories across usefulness, relevance, accuracy, and granularity to enable optimal data selection.

Table 3. Comparison of Data Evaluation and Filtering Strategies.

Strategy	Advantages	Limitations
Environment-based	Easy to implement with clear and well-defined reward signals.	Primarily focuses on final binary outcomes, neglecting process feedback and potentially accumulating errors.
Human or rule-based	Highly reliable with strong interpretability and customizable evaluation.	Relies on predefined metrics or human oversight, and requires complex rule design.
Model-based	Fully automated process without the need for human intervention.	Depends on model reliability and is affected by inherent biases within the model.

**Summary.** In the evaluation and filtering of trajectory data, different strategies offer distinct advantages and limitations, as summarized in Table 3. Environment-based evaluation methods are easy to implement, as they rely on clear binary feedback based on task success or failure. However, these methods focus primarily on final outcomes, and their discrete binary rewards limit feedback granularity, potentially masking reasoning errors in the trajectory generation process. A more robust approach is to develop process-based rewards that track decision-making steps and provide more informative evaluation signals. Human or rule-based methods provide adaptable and task-specific evaluation, offering high precision but requiring complex design and manual oversight. Model-based methods enable fully automated evaluation with rich interpretability, reducing human effort, yet their effectiveness is limited by the biases and errors inherent in the evaluating models. Each method presents trade-offs in scalability, accuracy, and implementation complexity, requiring careful selection and design based on the specific demands of the agent task. Additionally, integrating multiple approaches in a hybrid strategy can further enhance data quality and robustness, ensuring a more reliable evaluation framework tailored to different task needs.

**Low-quality Data Utilization.** Utilizing low-quality data has emerged as an effective strategy to augment training datasets. Early methods predominantly relied on expert-annotated or successful

trajectory data, resulting in underutilization of available trajectories. However, failed samples also contain latent instructive value, enabling agents to learn from both successful and erroneous cases, which enhances decision-making robustness and generalization capacity.

Some works create comparative datasets by pairing successful trajectories with failed ones, allowing agents to learn from both correct and incorrect behaviors. ENVISIONS [178] pair low-quality trajectories with correct solutions to create positive and negative sample pairs for training, enhancing the agent’s ability to distinguish between correct and incorrect actions. Similarly, AgentGym [170], IPR [177] and ETO [212] generate preference-based trajectory pairs and apply DPO to optimize agent policies. Other methods focus on correcting low-quality samples to improve training data quality. AGENTBANK [143] refines failed trajectories to regenerate interactions, thereby refining the dataset. Re-ReST [33] leverages a trained reflector model to iteratively correct low-quality data based on ground truth, environmental feedback, and task-specific information. ADASWITCH [146] also adopts an iterative correction process by employing a hybrid approach, where a small local model generates task steps, and a cloud-based large model corrects errors before reintegrating the refined data into the dataset. In addition to these correction methods, some approaches directly utilize error samples to teach models to recognize and understand failure scenarios. Agent-FLAN [22] incorporates frequently encountered negative samples into the context, explicitly teaching the model which actions should be avoided, thereby reducing hallucinations. NAT [158] leverages modified instruction prompts to generate incorrect responses, reinforcing the model’s ability to distinguish between valid and invalid outputs.

**Summary.** Leveraging low-quality data provides a way to enhance agent learning by extracting useful information from failure cases. Comparative approaches improve discrimination between correct and incorrect actions, correction-based methods refine training datasets for more accurate learning, and direct error utilization strategies reinforce robustness against failures. However, ensuring the consistency and reliability of self-generated failure cases, minimizing bias in correction processes, and balancing error utilization with preventing undesirable learning effects require careful and rigorous design.

**3.1.2 Trajectory-based Fine-Tuning for Agent.** Fine-tuning is a critical step in the optimization process for LLM-based agents, enabling open LLMs to adapt to specific agent tasks or data distributions. Most of the techniques used here are based on fine-tuning methods commonly used in LLMs, but are specifically applied to agent-related data (constructed in §3.1.1) to enhance the decision-making and task-solving capabilities of LLM-based agents.

**Mixing of General Instructions and Agent Trajectories.** Some studies [22, 216] have shown that fine-tuning LLM-based agents solely on agent-specific trajectory data may weaken the abilities of general language understanding and reasoning. To mitigate this issue, many works incorporate a mix of general instruction datasets and task-specific trajectory data during fine-tuning, which ensures that LLMs retain foundational capabilities while being optimized for agent-specific tasks. AgentTuning [199] combines AgentInstruct datasets with general-purpose instruction data to fine-tune LLaMA-2 models using instruction tuning, ensuring both task-specific alignment and foundational language abilities. AGENTBANK [143] integrates agent trajectory data, general instruction data, and code datasets to optimize LLaMA-2-Chat models for diverse tasks. Similarly, Zhou et al. [216] mix the agent trajectory datasets and general instruction datasets to fine-tune OpenChat and Llama-2 series models by LoRA. In summary, data mixing strategies enable fine-tuned LLM-based agents to preserve both foundational and complex agent capabilities, balancing general language understanding with task-specific optimization for enhanced versatility and robustness.

**The Usage of Fine-Tuning Techniques.** Based on studies on LLM-based agents, we categorized fine-tuning methods into three types: standard SFT, parameter-efficient fine-tuning (e.g., LoRA), and customized strategies tailored for specialized tasks, which will be discussed in detail.

**(1)Standard SFT.** The goal of SFT in LLM-based agents is to align pre-trained models with specific task requirements by minimizing the discrepancy between predicted and target outputs, typically in the form of instruction-output pairs or trajectory datasets. In this section, we define and classify standard SFT to include works that involve full parameter fine-tuning (Full Fine-Tuning) with high-quality training data and instruction tuning. We also classify works that directly mention using SFT-based fine-tuning without providing further technical details as part of the standard SFT category. Additionally, given that Behavior Cloning (BC) in imitation learning aligns with the supervised fine-tuning paradigm, we similarly include it within this category to provide a more comprehensive overview. Most works utilize standard SFT to fine-tune LLM-based agents using high-quality datasets. AgentTuning [199], Agent-FLAN [22] and AGENTBANK [143] fine-tune Llama-based models, relying on instruction data to align models with specific task requirements. Similarly, NAT [158], STE [154], and COEVOL [85] employ SFT to optimize models like Llama-2-Chat or Mistral-Instruct by leveraging trajectory datasets curated for task success and failure scenarios. In addition, AGILE [39], Retrospec [171], AgentGym [170], ETO [144] and ToRA [49] adopt BC, imitation learning-based SFT, to equip LLMs with foundational agent capabilities, enabling them to better handle agent-specific tasks. Standard SFT remains a widely adopted approach for fine-tuning LLM-based agents, leveraging diverse datasets such as instruction data, trajectory datasets, and behavioral examples. Its simplicity and effectiveness make it a foundational method in agent optimization, particularly for aligning models with task-specific requirements.

**(2)Parameter-efficient fine-tuning.** PEFT methods, such as LoRA and QLoRA, optimize only a small subset of parameters while keeping the majority of the model frozen. These techniques significantly reduce computational and memory costs, making them ideal for fine-tuning large-scale LLMs. By focusing on specific task-relevant parameters, parameter-efficient methods achieve comparable performance to full-parameter optimization while improving efficiency. Several works employ PEFT to optimize LLM-based agents. For example, SMART [197], FireAct [14], Re-ReST [33] and Agent Lumos [192] utilize LoRA to fine-tune models from the Llama-2 and LLaMA-3 series, OpenChat-v3.2, and CodeLlama. Various configurations are set by these works, including chat and instruct versions, and different parameter sizes, effectively aligning the models with task-specific requirements. Additionally, AgentOhana [202] applies QLoRA to optimize xLAM-v0.1 models, demonstrating its ability to fine-tune large models efficiently under resource constraints. In summary, PEFT offers a practical solution for optimizing LLM-based agents, balancing performance and efficiency in resource-constrained scenarios.

**(3)Customized fine-tuning.** Customized fine-tuning methods are tailored to specific tasks or ideas, incorporating unique strategies or objective functions to better align models with agent-specific requirements. These approaches often include task-specific modifications, such as combining trajectory data with instruction datasets or introducing additional constraints like regularization terms. ATM [219] designs Multi-agent Iterative Tuning Optimization (MITO) loss, which combines standard SFT with KL regularization to balance task-specific optimization and generalization, to fine-tune Llama-2-7B. Leveraging contrastive learning and RL-free optimization, ENVISIONS [178] achieves efficient iterative refinement of Llama-2 through cyclic fine-tuning on dynamically updated trajectory pairs. These customized fine-tuning strategies highlight the flexibility of adapting LLMs to agent-specific tasks by integrating task-relevant objectives and iterative optimization processes.

**Summary.** Conventional fine-tuning-based optimization has effectively enhanced LLM-based agents by leveraging high-quality trajectory data, different training methodologies, and efficient parameter adjustments. This approach ensures strong task alignment, controlled optimization, and

adaptability to specific objectives, making it a reliable method for refining agent behavior. However, it also face inherent limitations that affect their effectiveness and adaptability. These methods rely heavily on high-quality trajectory data, making performance contingent on data availability and curation. Challenges such as error accumulation, overfitting to specific datasets, and difficulty in adapting to dynamic environments can hinder generalization across diverse tasks. Additionally, fine-tuning primarily aligns models with static objectives and lacks interactive feedback mechanisms, limiting the agent’s ability to refine its behavior based on real-time task execution. As LLM-based agents are expected to handle increasingly complex and evolving tasks, exploring more adaptive optimization strategies becomes essential to further enhance their decision-making capabilities.

### 3.2 Reinforcement Learning-based Optimization

To mitigate the limitations of conventional fine-tuning, RL-based optimization methods have emerged as a promising approach, which enable LLM-based agents to learn directly from interactions with environment/human feedback, and leverage rewards and penalties to dynamically refine behavior. Unlike static fine-tuning, RL techniques encourage exploratory learning, allowing agents to discover novel policies and adapt to unseen tasks or dynamic conditions. By aligning model outputs with explicit rewards or human-aligned preferences, they not only enhance task-specific performance but also foster greater flexibility and robustness in complex, evolving environments.

We categorize RL-based optimization methods into two primary approaches: **Reward-function-based optimization** and **Preference-alignment-based optimization**. The former leverages traditional reinforcement learning techniques, such as PPO, to guide iterative learning process through explicitly defined reward signals. In contrast, the latter focuses on newer methods like DPO, relying on preference datasets to align model outputs with human preferences, bypassing the need for traditional reward modeling. The workflow of both approaches is illustrated in Figure 3.

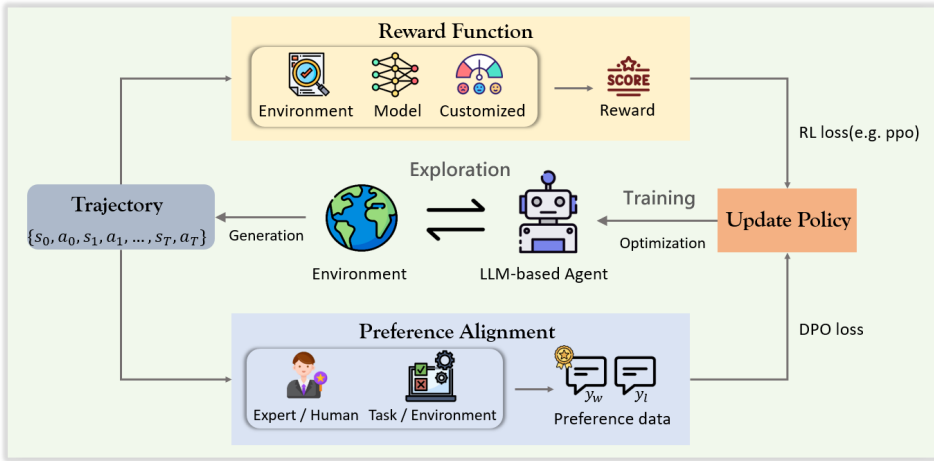


Fig. 3. Workflow of Reinforcement Learning-based Optimization for LLM-based Agents.

**3.2.1 Reward Function-based Optimization.** Reward function-based optimization methods leverage explicit reward signals to refine the behavior of LLM-based agents, enabling their adaptation to complex tasks and dynamic environments. Drawing from traditional RL paradigms, these methods employ algorithms like PPO or Actor-Critic to iteratively optimize agents’ policies and

adjust LLMs’ parameters. By framing the LLM itself as the agent, these approaches utilize diverse reward sources, including environmental feedback, model-generated signals, and customized reward functions, as summarized in Table 4. In this part, we summarize representative works categorized by the nature of their reward signals and construction.

Table 4. Summary of Reward Function-based Optimization Methods for LLM-based Agents. Note: MA-Multi-Agent Framework.

Methods	Reward Source	RL Algorithm	MA	Base Model
CMAT [92]	Environment-based	Actor-Critic	✓	Tinyagent-1.8B/7B
StepAgent [29]	Model-based	DPO+PPO	/	Mistral-7B, Llama3-8B
WebRL [126]	Model-based	Actor-Critic	✓	GLM-4-9B, Llama-3.1-8b/70B
CORY [111]	Customized, Model-based	PPO	✓	GPT 2-Large, Llama-2-7B-Chat
SaySelf [179]	Customized	PPO	✓	Mistral-7B, Llama-3-8B
AgentGym [170]	Environment-based	AgentEvol	/	Llama-2-7B-Chat
GELI [78]	Model-based	PPO	/	Llama-2
AGILE [39]	Customized	PPO	/	Vicuna-13B, Meerkat-7B

**Environment-based rewards.** One of the most common approaches leverages environmental feedback as the primary reward source. For example, CMAT [92] employs an actor-critic framework within a multi-agent collaboration setting, where agents adaptively interact with their environment and receive rewards derived from task performance and interaction outcomes.

Retrospec [171] employs Implicit Q-Learning in an offline RL framework, mitigating Q-function overestimation through temporal difference (TD) error minimization on fixed environmental trajectories. This approach stabilizes policy updates by constraining value estimation within empirically observed state-action pairs. AgentGym [170] optimizes itself using the AgentEvol algorithm, which includes an optimization objective function based on environment trajectory rewards. Furthermore, AGILE [39] combines action correctness with designed reward for expert assistance requests. Using these explicit feedback signals in conjunction with PPO, AGILE optimizes the ability of agent to solve problems independently and seek the help of external experts in time. Environment-based rewards are straightforward to obtain as they directly rely on task feedback from the environment. However, these rewards are often discrete and outcome-focused, which may limit their ability to provide detailed guidance during intermediate steps of complex tasks.

**Model-based rewards.** Model-based rewards derive their signals from LLMs or models trained to evaluate agent behaviors, providing implicit or explicit feedback to guide policy optimization. These methods are particularly useful in scenarios where direct environmental rewards are sparse or unavailable. StepAgent [29] utilizes an inverse reinforcement learning (IRL) framework, where a discriminator predicts the likelihood of an agent’s trajectory aligning with expert behavior, serving as a reward signal. It employs step-wise optimization by comparing expert and agent actions to generate fine-grained intermediate rewards, enabling step-level adjustments. Besides, Jensen–Shannon (JS) divergence ensures that the agent’s behavior remains close to expert trajectories, while the discriminator distinguishes between them, forming a structure similar to adversarial networks. In addition, WebRL [126] trains a self-supervised Optimal Reward Model (ORM) to evaluate trajectory quality and generate implicit feedback. Using an actor-critic framework, the model updates policy weights based on ORM-generated rewards, while integrating Kullback-Leibler(KL) divergence constraints to maintain policy stability across iterative updates. In summary, model-based rewards offer flexibility by decoupling reward generation from direct environmental feedback, allowing for task-specific evaluations. However, their effectiveness depends on the robustness of the reward models, which require careful design and optimization.

**Customized reward functions.** Customized reward functions are tailored to specific tasks, allowing researchers to design multi-objective signals that optimize agent behavior across multiple dimensions. These functions often go beyond task success rates to incorporate additional objectives, such as policy stability, collaboration efficiency, and task-specific performance. GELI [78] combines Global Explicit rewards (GE), such as session-level user satisfaction, with Local Implicit feedback (LI) derived from multimodal signals. This multi-objective design decomposes rewards into fine-grained components, enabling more effective optimization of long-term conversational goals through PPO. Similarly, AGILE [39] uses a specialized reward function to guide agents in seeking expert assistance. By integrating binary task success indicators and a cost-based penalty for expert queries, AGILE incentivizes efficient problem-solving with minimal resource usage. CORY [111] designs shared rewards in a multi-agent framework, where a leader and observer alternate roles to optimize behavior. KL divergence constraints are used to minimize policy deviations, ensuring stability during updates. Although CORY employs PPO, the framework is compatible with other RL algorithms, enhancing its generalizability. To refine the confidence estimation in generated outputs, SaySelf [179] introduces a reward function emphasizing answer accuracy and penalizing overconfidence in erroneous predictions. Optimized through PPO, the design reduces misleading outputs, and ensures robustness in decision-making under uncertainty. These customized reward function methods demonstrate the adaptability of RL-based frameworks to complex, real-world challenges, showcasing innovative approaches to reward design that cater to diverse agent optimization needs.

**Summary.** Reward function-based optimization methods leverage traditional RL algorithms to iteratively refine LLM-based agent behavior using explicit reward signals. These methods are effective in adapting agents to complex tasks and dynamic environments. However, designing appropriate reward functions requires careful consideration of the target task to ensure that the reward structure captures all relevant aspects, avoiding overly narrow or generic criteria. Despite their effectiveness, these methods face challenges such as the extensive interaction data and computational resources required by algorithms like PPO, raising scalability concerns for large-scale LLM agents. Additionally, reliance on auxiliary models, such as value functions in Q-Learning, adds complexity, requiring robust model integration and fine-tuning. These challenges highlight the need for advancements in both algorithmic efficiency and reward signal design.

**3.2.2 Preference Alignment-based Optimization.** Preference alignment-based optimization provides an alternative paradigm of traditional RL by directly aligning model behavior with human or expert preferences, without relying on explicit reward signals. Inspired by RLHF, preference alignment-based optimization borrows key elements such as iterative policy updates and feedback-driven learning. By utilizing preference data constructed through pairwise comparisons or other ranking mechanisms, these methods optimize policies to better reflect desired outcomes, simplifying the process by directly optimizing preference datasets through techniques like DPO. DPO can be viewed as a reinforcement learning variant tailored for offline datasets, where preferences replace traditional reward signals to guide policy refinement. This offline optimization approach enables effective agent fine-tuning without requiring online deployment, making it particularly suitable for large-scale LLM-based agents. The optimization process can be divided into two main stages: **Preference Data Construction** and **Policy Optimization** as follows.

- **Preference Data Construction.** The first step is the construction of preference data, which captures user or expert preferences in the form of pairwise comparisons. These comparisons are typically represented as tuples  $(x, y_w, y_l)$  where  $x$  is the input instruction/question,  $y_w$  denotes the preferred or better(win) response, and  $y_l$  represents the less preferred or suboptimal(lose) response. Some studies alternatively use  $a^+$  and  $a^-$  to denote the preferred and less preferred outputs, respectively. This structured format ranks outputs explicitly,

providing a clear foundation for preference-based optimization. These pairs are generated from a variety of sources, including expert trajectories, human feedback, task-specific metrics, or agent-generated exploration trajectories. For the exploration stage, methods like Monte Carlo Tree Search (MCTS) are often employed to systematically generate diverse and high-quality trajectories. These trajectories are then evaluated based on task success rates, environment feedback, or other metrics to construct  $y_w$  and  $y_l$ . User preferences collected through direct comparisons or rankings further enhance the alignment of agents' actions with human expectations.

- **Policy Optimization via DPO.** Following the construction of preference data, the next step is optimizing the agent's policy to align with these preferences. DPO serves as a widely adopted approach for this purpose, leveraging preference pairs to refine LLMs parameters. By iteratively optimizing the policy, DPO enables LLM-based agents to prioritize desired outputs and adhere to task-specific requirements. The optimization objective of DPO is formalized as follows:

$$\mathcal{L}_{DPO}(\pi_{\theta}; \pi_{ref}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y_w|x)}{\pi_{ref}(y_w|x)} - \beta \log \frac{\pi_{\theta}(y_l|x)}{\pi_{ref}(y_l|x)} \right) \right] \quad (1)$$

where  $\mathcal{D}$  is the preference dataset,  $\pi_{\theta}(y|x)$  denotes the probability assigned by the current policy to output  $y$  given input  $x$ , and  $\pi_{ref}(y|x)$  serves as a reference from the LLM-based agent. The regularization term  $\beta$  ensures the updates remain consistent with the original policy, while the sigmoid function  $\sigma$  provides stability during training. DPO bypasses the need for additional interactions with the environment through directly preference-based policy optimization. Its regularization mechanism ensures stability and avoids significant deviations from the reference model.

We discuss the representative works under this paradigm, categorized by the criteria used for constructing preference datasets, including those based on expert, human and task feedback, as summarized in Table 5.

Table 5. Summary of Preference Alignment-based Optimization Methods for LLM-based Agents. Note: MA-Multi-Agent Framework.

Methods	Preference Criteria	Algorithm	MA	Base Model
StepAgent [29]	Expert	DPO+PPO	/	Mistral-7B, Llama3-8B
AgentQ [125]	Environment	DPO	/	Llama-3 70B, xLAM-v0.1-r
OPTIMA [19]	Task, Human	DPO	✓	Llama-3-8B
IPR [177]	Expert, Environment	DPO	/	llama-2-7B
Re-ReST [33]	Environment	DPO	✓	Llama-2-7/13B, Llama-3-8B, Codellama-13B, VPGen
DMPO [138]	Expert	DMPO	/	Llama-2-7B-Chat, Mistral-7B-Instruct-v0.2
EPO [212]	Task, Environment	DPO	✓	Llama-2-7B
ATM [219]	Human	DPO	✓	Mistral-7B-Chat
AMOR [51]	Expert, Human	KTO	✓	Llama-2-7B/13B-Chat

**Optimization Based on Expert or Human Preferences.** Preference data in this category rely on expert trajectories or human-defined preferences to construct positive and negative samples, providing a clear benchmark for optimization. Positive samples  $y_w$  are typically derived directly from expert behavior or human annotations, representing ideal responses or actions. Negative samples  $y_l$  are often derived from agent-generated error or suboptimal trajectories or noisy data. DMPO [138] refines RL objectives by replacing policy constraints with State-Action Occupancy Measure (SAOM) constraints, using expert trajectories to define preferred and dispreferred pairs.



SAOM constraints guide the model to mimic expert state-action distributions, especially in unexplored states, ensuring actions align with expert trajectories. By introducing a simplified loss function with length normalization, DMPO directly maximizes the likelihood of preferred trajectories, reducing errors in multi-turn tasks. IPR [177] constructs preference data by comparing agent actions with expert trajectories at each step. If the agent’s action yields a lower reward than the expert’s, it is marked as an error, forming a preference pair with the expert’s action as the positive sample. These step-level preferences are optimized using both outcome-DPO loss and step-DPO loss, ensuring fine-grained alignment with expert standards. Distinctively, AMOR [51] adopts a two-stage approach for preference alignment, combining SFT during pre-training with Kahneman-Tversky Optimization (KTO) in adaptation. In the adaptive stage, human-annotated binary feedback (right/wrong) forms the basis for optimizing module-specific parameters. KTO improves sensitivity to negative samples and directly maximizes human-defined utility, diverging from traditional preference-based log-likelihood methods.

While preference data construction based on expert and human preferences leverages high-quality trajectories to ensure accurate optimization, it often suffers from limitations in scalability and coverage, as expert trajectories may not sufficiently represent diverse scenarios or unexplored states in complex environments.

**Optimization Based on Task or Environment Preferences.** This category involves constructing preference datasets using task-specific metrics or feedback derived directly from interactions with the environment. These approaches leverage rewards, success rates, or other performance indicators to evaluate trajectories, creating preference pairs for training through DPO algorithms. By incorporating environment-specific signals, these methods ensure adaptability to dynamic and task-oriented contexts. Some works [33, 144, 170] construct preference data using straightforward approaches based on task success or environmental feedback, using DPO to optimize agent policies in alignment with task objectives. Environment Preference Optimization (EPO) [212] uses environment-based feedback, ranking outputs through a reward model to identify winning and losing samples. It combines DPO with a token-level alignment loss as to ensure stable training and precise task alignment. Other approaches introduce more structured methods for preference construction. For example, AgentQ [125] employs MCTS to explore various branches of action trajectories, generating preferences by combining AI process feedback and success rates. Similarly, OPTIMA [19] generates trajectory pairs by analyzing MCTS trees, ranking node pairs based on criteria such as shared ancestry and significant reward differences. The top-ranked pairs are selected for DPO training to improve efficiency and task performance. ATM [219] uses perplexity scores as task reward to distinguish between high- and low-quality outputs, forming binary preference pairs that are optimized using DPO to enhance reasoning and task adaptability.

**Summary.** RL-based optimization methods have demonstrated significant potential in improving the decision-making capabilities of LLM-based agents. By incorporating feedback signals, such as reward functions or preference alignment, these methods enable agents to adapt to complex and dynamic tasks. Worth noting is that the choice of these method depends on task complexity and data availability: reward-function-based methods excel in handling complex scenarios but require extensive data and computational resources, while preference-alignment approaches like DPO simplify training and align outputs with human preferences but rely heavily on the quality and coverage of preference data, limiting adaptability in highly dynamic tasks.

### 3.3 Hybrid Fine-Tuning Optimization

While traditional SFT provides a stable method for initializing LLM-based agents, it often struggles with dynamic tasks requiring complex decision-making. On the other hand, RL excels in exploring complex scenarios but typically demands extensive data and computational resources, making it

less suitable for direct initialization. To overcome these limitations, hybrid fine-tuning strategies combine the strengths of both SFT and RL, creating a more flexible and effective framework. This section introduces representative works that exemplify these hybrid approaches.

Most works in hybrid fine-tuning follow a sequential approach, starting with a warm-up stage where behavior cloning-based SFT is applied to equip LLMs with foundational capabilities using high-quality datasets, such as expert trajectories. In the subsequent RL stage, algorithms like PPO or DPO refine the agent’s policy for task-specific objectives or dynamic environments. This paradigm, similar to Reinforcement Fine-Tuning (RFT) introduced by OpenAI, has gained widespread popularity for combining SFT and RL to enhance LLMs’ capabilities. For example, ReFT [152], AgentGym [170], ETO [144], Re-ReST [33], AGILE [39], and AMOR [51] all utilize SFT during the warm-up phase to train models on expert trajectories or curated datasets. Following initialization, these methods apply various RL strategies. ETO and Re-ReST use DPO for preference alignment, AGILE employs PPO for decision-making, and AMOR applies Kahneman-Tversky Optimization (KTO) with binary feedback to align outputs with human preferences.

Some studies adopt iterative approaches to refine the hybrid fine-tuning paradigm, alternating between SFT and RL phases to enhance performance. OPTIMA [19] exemplifies this by combining Iterative SFT and Iterative DPO, enabling the LLM agent to learn from optimal trajectories through SFT while refining its understanding using DPO based on comparative preferences. Similarly, IPR [177] incorporates step-level rewards, starting with ReAct-style expert trajectories for SFT and iteratively refining the policy through Outcome-DPO and Step-DPO losses. These iterative cycles ensure continuous alignment with task objectives by addressing both global and fine-grained improvements.

Additionally, some works incorporate specialized mechanisms within the hybrid fine-tuning framework to address specific challenges. For example, Retrospect [171] starts with SFT on high-quality expert trajectories to establish foundational knowledge and task-specific skills. During the post-training phase, it integrates offline RL by training an action-value function through IQ-learning, allowing the model to combine fixed experience data with LLM outputs via an RL Critic for weighted scoring, refining decision-making in dynamic environments. ENVISIONS [178] uses SFT to optimize self-exploration-derived solutions, selecting positive and negative pairs from a candidate pool based on task success metrics. To simplify traditional RL, it replaces reinforcement learning with an RL-Free loss combined with an MLE-based fine-tuning loss for the final optimization.

In conclusion, hybrid fine-tuning strategies combine the strengths of SFT and RL, allowing LLM-based agents to balance structured guidance with adaptive optimization. While these methods offer greater flexibility for complex, dynamic tasks, they also face challenges such as increased computational costs and reliance on high-quality preference or reward data. Future research should focus on developing more effective methods to combine fine-tuning and reinforcement learning. This hybrid paradigm holds significant potential for advancing agent optimization by enhancing the versatility and performance of LLM-based systems. The effectiveness of these strategies depends largely on their design and integration, making this an essential direction for future innovations.

#### 4 Parameter-Free Optimization of LLM-based Agents

Besides parameter-driven optimization methods, which are critical for evolving LLM-based agents, **Parameter-Free Optimization** offers a promising alternative by optimizing the agent’s behavior without modifying its underlying model parameters. This method improves performance by adjusting the model’s inputs, context, and task interactions, primarily through natural language prompts, making it especially efficient in resource-constrained environments. This section examines various approaches to Parameter-free optimization, which focus on enhancing agent capabilities through prompt engineering and incorporating various types of information within the prompts, such as

Table 6. Comparison of Parameter-Free Optimization Methods.

Optimization Strategy		Description	Example Works
Experience-based		Learning from past interactions and historical experience	Optimus-1 [90], Agent Hospital [82], ExpeL [211], AutoManual [16], AutoGuide [42], Experiential Co-Learning [127]
Feedback-based	Self-Reflection	Self-evaluating and correcting to refine current behavior	Reflexion [139], QueryAgent [62], SAGE [93], Agent-pro [207], Recon [159], Symbolic Learning [218], NLRL [40]
	External Feedback	Modifying behavior via external signals and critiques	Retroformer [190], COPPER [10], InteRecAgent [63], CoE [175]
	Meta-Prompt	Iteratively refining global instructions/prompts	MetaReflection [53], Retroformer [190], OPRO [181], Symbolic Learning [218]
Tool-based		Using external tools and optimizing selection strategies	TPTU [135], AVATAR [168], Lyra [213], Middleware [50], AgentOptimizer [206], VideoAgent [36], AutoAct [130]
RAG-based		Retrieving and incorporating external knowledge	AutoRAG [72], Self-RAG [8], RaDA [73], Rap [71], PaperQA [76], MALADE [23], EMG-RAG [162]
Multi-Agent collaboration		Multiple role-specific agents communicate and collaborate	MetaGPT [57], ChatDev [128], MapCoder [66], DyLAN [107], MacNet [129], Agentverse [18], CAPO [100], SMOA [80], MAD [91], AutoGen [167], Multi-AI Agent [198]

feedback, tools, external knowledge and agent collaboration. Unlike other surveys [64, 169, 210] that categorize methods based on agent structure or architecture, we categorize the methods based on the strategies and techniques employed to optimize the agent’s behavior, including experience-based optimization, feedback-based optimization, tool-based optimization, RAG-based optimization, and multi-agent collaborative optimization.

#### 4.1 Experience-based Optimization

Experience-based optimization leverages historical data, trajectories, or accumulated knowledge to improve LLM-based agents by deriving insights from past experiences. By storing and analyzing successes and failures, often through memory modules, agents can refine strategies, enhance long-term decision-making, and adapt to evolving tasks. This method provides a powerful mechanism for agents to generalize across tasks and domains.

Optimus-1 [90] utilizes a multimodal memory module to convert exploration trajectories into hierarchical knowledge graphs, which assist in task planning and prompt generation for the agent. Similarly, Agent Hospital [82] integrates a medical records library and an experience repository to refine guidelines based on successful and failed cases, optimizing agent decision-making in medical scenarios. To extract actionable insights, ExpeL [211] autonomously collects knowledge from training tasks, and recalls these insights during inference. AutoManual [16] organizes task rules into manuals that serve as prompts to guide the agent’s future actions and update system regulations. AutoGuide [42] uses offline trajectory data to map contexts to optimal actions, allowing agents to retrieve top-ranked guidelines during inference and adapt dynamically. In addition, Experiential Co-Learning [127] introduces a framework in which instructor and assistant agents collaboratively gather shortcut-oriented experiences from their historical trajectories. These experiences are reused to inform future task execution, enabling agents to iteratively improve through shared and accumulated knowledge.

#### 4.2 Feedback-based Optimization

Feedback-based optimization enhances LLM-based agents by leveraging feedback for self-reflection, correction, and iterative improvement. Through dynamic updates, agents evolve their policies and abilities to adapt to complex tasks. These approaches are close to parameter-driven iterative optimization methods discussed in section 3.1, but instead of modifying LLM parameters, it utilizes natural language-based feedback and reflection to iteratively guide and optimize agent behavior. We categorize related works into three types: feedback-based self-reflection optimization, external

feedback-based optimization, and meta-prompt optimization, based on the strategies used to optimize agent behavior and how feedback triggers the next round of improvement.

**4.2.1 Feedback-based Self-Reflection Optimization.** Feedback-based self-reflection optimization focuses on the agent’s ability to reflect on its actions and decisions using feedback from the environment or its own evaluations. The agent uses this feedback to identify areas for improvement and adjust its behavior through self-correction and evolution, enhancing its adaptability and performance in dynamic environments.

Reflexion [139] and QueryAgent [62] use environmental and task-specific feedback to drive self-reflective adjustments. By converting task outcomes or heuristic evaluations into textual corrections, these methods integrate feedback into decision-making, improving adaptability and contextual understanding. In SAGE [93], the checker agent provides iterative feedback on the current solution, while the assistant agent generates self-reflection for further feedback to support self-evolution. Other works expand self-reflection by incorporating higher-order reasoning and strategy-level optimization. Agent-pro [207] iteratively reflects on past trajectories and beliefs, allowing the agent to adjust irrational strategies and refine decision-making. Similarly, Recon [159] introduces first- and second-order perspective shifts, where agents evaluate their reasoning and consider external viewpoints for strategic adjustments. Additionally, methods like Symbolic Learning [218] and NLRL [40] employ natural language prompts to mimic parameter optimization, such as gradient updates and reward functions. These methods use carefully designed prompts to evaluate the current decision-making process and treat the resulting evaluations as feedback, enabling agents to iteratively refine their strategies and improve task performance.

**4.2.2 External Feedback-based Optimization.** External feedback-based optimization leverages evaluative signals from external models, agents, or frameworks to refine the behavior of the actor agent. Inspired by the actor-critic paradigm in RL, these methods integrate external reflections and corrections to enhance the robustness and adaptability of LLM-based agents, which are commonly employed in multi-agent scenarios and represent frequent technologies within multi-agent collaborative optimization frameworks, which will be discussed in §4.5.

Retroformer [190] utilizes a retrospection model trained with PPO to analyze failures and provide refined feedback, enabling the primary agent to iteratively adjust its actions based on retrospective insights. Similarly, COPPER [10] employs a shared reflector module to generate counterfactual feedback, refine prompts, and store improvements in a memory module, enhancing the actor agent’s task performance and reliability. InteRecAgent [63] involve a critic LLM overseeing the actor agent’s behavior to prevent instruction violations or incorrect tool usage, thereby enhancing robustness and error correction. In CoE [175], a multi-agent framework is proposed to address complex Operations Research, where a conductor orchestrates task assignment among specialized agents, and an evaluator provides feedback to iteratively refine the outputs through backward optimization.

**4.2.3 Meta-Prompt Optimization.** Meta-prompt optimization focuses on refining global instructions or meta-prompts to enhance the generalization capabilities of LLM-based agents. By iteratively adjusting prompts based on feedback, it enables agents to adapt more effectively to diverse and dynamic tasks, focusing on optimizing their behavior from a broader and more general perspective.

MetaReflection [53] extracts information from failed trials to create optimized prompts, iteratively incorporating them to enhance task performance. Retroformer [190] refines meta-prompts by training a retrospective agent with PPO, using failure analyzes to generate improved instructions that guide agents in avoiding repeated mistakes. Other approaches focus on simulating optimization

processes through meta-prompts. For instance, OPRO [181] refines prompts by analyzing task accuracy and generating instructions tailored to improve outcomes iteratively. Symbolic Learning [218] uses natural language to simulate optimization concepts such as gradient updates and loss functions, enabling agents to iteratively adjust their strategies through language-driven prompts.

### 4.3 Tool-based Optimization

One of the key distinctions between LLMs and LLM-based agents is their ability to utilize external tools. These tools, including calculators, search engines and code interpreters (e.g., Python), APIs, and domain-specific modules, enable agents to perform tasks requiring external computation, dynamic information retrieval, or specialized functionalities. Applying these tools, agents can significantly enhance their problem-solving scope and capabilities, making tool utilization a cornerstone of agent optimization. In this section, we introduce works that enhance agents' performance by optimizing their tool usage and selection strategies.

TPTU [135] enhances task planning and tool usage by optimizing task decomposition and tool invocation. AVATAR [168], Lyra [213], and Middleware [50] refine tool strategies to address existing limitations. AVATAR employs a comparator to analyze performance differences between sample pairs, attributing discrepancies to tool usage issues and providing actionable improvements. Lyra focuses on correcting tool errors in formal proofs using predefined prover tools like Sledgehammer. Middleware introduces error feedback mechanisms to align tool inputs and outputs across steps, mitigating execution errors and improving system robustness. Additionally, AgentOptimizer [206] optimizes tool usage by treating functions as learnable weights, enabling agents to iteratively refine function sets based on execution history and task performance without modifying core LLM parameters. Other methods prioritize efficient tool integration in dynamic environments. VideoAgent [36] employs a minimal yet sufficient tool set, incorporating specialized models for visual question answering and object memory queries. AutoAct [130] automates task trajectory generation using a pre-assembled library of tools, such as search engines and code interpreters.

### 4.4 RAG-based Optimization

RAG has emerged as a powerful paradigm for enhancing LLM-based agents by integrating retrieval with generative processes. By dynamically incorporating external knowledge from large-scale databases or domain-specific corpora, RAG overcomes the limitations of fixed pre-trained knowledge, improving adaptability in evolving environments. These approaches build on in-context learning, augmenting decision-making and reasoning in knowledge-intensive tasks. In this section, we explore representative works that apply RAG to optimize agent performance, focusing on strategies for integrating retrieval mechanisms with prompt-based enhancements.

Some methods focus on optimizing retrieval configurations and augmenting language model outputs. AutoRAG [72] automates the selection of RAG modules, systematically evaluating combinations of retrieval techniques, re-ranking strategies, and expansion methods to identify optimal configurations. Self-RAG [8] integrates retrieval with self-reflection, allowing agent to adaptively refine content through iterative feedback. Other frameworks emphasize retrieval-enhanced task planning and execution. RaDA [73] and Rap [71] utilize past experiences and dynamic retrieval to decompose tasks and generate contextually informed actions, enhancing agent capabilities in iterative planning scenarios. Similarly, PaperQA [76] extends RAG applications to scientific domains by performing retrieval across full-text articles, ensuring precise and contextually relevant answers to scientific queries. In multi-agent contexts, MALADE [23] leverages retrieval to enhance coordination among agents in domain-specific problems, such as adverse drug event identification. By extracting relevant information from external data sources, agents collaboratively address complex

tasks. Similarly, EMG-RAG [162] employs a RL agent to optimize memory selection, enhancing contextual prompts for agent interaction.

#### 4.5 Multi-Agent Collaborative Optimization

Multi-agent frameworks are widely used for handling complex, dynamic tasks through collaboration. By distributing responsibilities across specialized agents, they enable parallel processing, information sharing, and adaptive role assignment, reducing cognitive and computational burdens. In LLM-based multi-agent optimization, prompts define agent roles, coordinate interactions, and establish shared goals. Through iterative collaboration and feedback, these frameworks enhance decision-making, enabling collective performance beyond individual capabilities.

Multi-agent collaborative optimization enhances agent performance in complex tasks by leveraging distributed roles and iterative interactions. In programming and software development, MetaGPT [57], ChatDev [128], and MapCoder [66] assign specialized roles to streamline workflows. These methods decompose tasks into modular phases, integrating structured prompts, role-based collaboration, and iterative refinement to simulate human-like software development processes, including retrieval, planning, code generation, debugging, and documentation. To optimize agent organization and collaboration, DyLAN [107] and MacNet [129] dynamically structure agent networks to identify the most effective combinations of agents. DyLAN introduces agent importance scores and early stopping mechanisms, allowing for efficient selection of key contributors, while MACNET organizes agents using directed acyclic graphs, streamlining interactions for more effective reasoning. Besides, Agentverse [18] and CAPO [100] focus on planning and task execution, allowing agents to iteratively evaluate and improve solutions. SMoA [80] and MAD [91] refine multi-agent interaction by introducing specialized frameworks for decision-making. MAD establishes a debate framework where multiple agents present arguments, guided by a judge agent who manages the debate and selects the optimal solution. Notably, the MAD framework has been frequently adopted in subsequent works [13, 88], demonstrating its effectiveness in multi-agent optimization tasks.

General-purpose frameworks like AutoGen [167] and Multi-AI Agent [198] design adaptive and flexible multi-agent systems for diverse tasks, ranging from programming to question answering. By dynamically coordinating specialized agents, these frameworks demonstrate the potential of multi-agent collaboration to address complex problems while achieving efficient and scalable optimization.

### 5 Datasets and Benchmarks

Datasets and benchmarks are essential for evaluating and enhancing LLM-based agents. In this section, we summarize the commonly used datasets and benchmarks from two perspectives: (1) datasets and benchmarks for evaluation, which assess agent performance across various tasks, (2) datasets for agent tuning, which are used for fine-tuning and optimizing the agent specific abilities in relevant tasks and environments.

#### 5.1 Datasets and Benchmarks for Evaluation

**5.1.1 General Evaluation Tasks.** In Table 7, we summarize common evaluation datasets, categorized by general task domains such as mathematical reasoning, Question Answering (QA) tasks, multimodal tasks, programming, etc. Additionally, following AgentBank [143], tasks are classified based on action space, distinguishing between tasks with a continuous action space (including natural language and code) and those with a predefined discrete action space.

**Math:** The mathematical reasoning datasets assess agents’ ability to perform multi-step reasoning and solve complex problems. GSM8K [26] offers grade school-level arithmetic problems, testing

Table 7. Summary of Commonly Used Datasets and Environments for Evaluation Tasks.

Domain	Datasets/Environments	Task Type	Action Space	Data Link
Math	GSM8K [26]	Mathematical reasoning	Continuous	<a href="#">Link</a>
	AsDIV [116]	Mathematical reasoning	Continuous	<a href="#">Link</a>
	SVAMP [123]	Mathematical reasoning	Continuous	<a href="#">Link</a>
	MATH [55]	Mathematical reasoning	Continuous	<a href="#">Link</a>
	AIME	Mathematical reasoning	Continuous	<a href="#">Link</a>
QA	HotpotQA [187]	Open-domain QA	Continuous	<a href="#">Link</a>
	StrategyQA [46]	Open-domain QA	Continuous	<a href="#">Link</a>
	MMLU [54]	Multiple choice QA	Continuous	<a href="#">Link</a>
	TruthfulQA [96]	Multiple choice QA, Open-domain QA	Continuous	<a href="#">Link</a>
	TriviaQA [70]	Open-domain QA	Continuous	<a href="#">Link</a>
	PubMedQA [68]	Domain-Specific QA	Continuous	<a href="#">Link</a>
	MuSiQue [151]	Commonsense Reasoning	Continuous	<a href="#">Link</a>
	2WikiMultihopQA [56]	Open-domain QA	Continuous	<a href="#">Link</a>
	QASPER [27]	Domain-Specific QA	Continuous	<a href="#">Link</a>
Code	ARC [24]	Commonsense Reasoning	Continuous	<a href="#">Link</a>
	SWE-bench [184]	Code Generation, Bug Fixing	Continuous	<a href="#">Link</a>
	HumanEval [17]	Code Generation	Continuous	<a href="#">Link</a>
	LiveCodeBench [118]	Code Execution, Bug Fixing	Continuous	<a href="#">Link</a>
	BIRD [81]	Text-to-SQL Conversion	Continuous	<a href="#">Link</a>
Tool-use	InterCodeSQL [185]	Code Generation	Continuous	<a href="#">Link</a>
	T-Eval [21]	Tool-use	Continuous	<a href="#">Link</a>
	MINT-Bench [160]	Tool-use	Continuous	<a href="#">Link</a>
	ToolEval [131]	Tool-use	Continuous	<a href="#">Link</a>
	API-Bank [84]	Tool-use	Continuous	<a href="#">Link</a>
Web	WebShop [188]	Web Navigation	Discrete	<a href="#">Link</a>
	WebArena [217]	Web Navigation	Discrete	<a href="#">Link</a>
	Mind2Web [28]	Web Navigation	Discrete	<a href="#">Link</a>
	MiniWoB++ [97]	Web Navigation	Discrete	<a href="#">Link</a>
Environment Interaction	ScienceWorld [157]	Scientific Experiment	Discrete	<a href="#">Link</a>
	ALFWorld [141]	Embodied AI	Discrete	<a href="#">Link</a>
	TDW-MAT [201]	Embodied AI	Discrete	<a href="#">Link</a>
	C-WAH [201]	Embodied AI	Discrete	<a href="#">Link</a>
	ALFRED [140]	Embodied AI	Discrete	<a href="#">Link</a>
	RLCard (Blackjack) [200]	Game	Discrete	<a href="#">Link</a>
	RLCard (Limit Texas Hold'em) [200]	Game	Discrete	<a href="#">Link</a>
	Breakthrough [77]	Game	Discrete	<a href="#">Link</a>
Multimodal	VQA v2.0 [7]	Visual QA	Continuous	<a href="#">Link</a>
	A-OKVQA [137]	Visual QA	Continuous	<a href="#">Link</a>
	ScienceQA (IMG) [108]	Visual QA	Continuous	<a href="#">Link</a>
	EgoSchema [114]	Video Understanding, Visual QA	Continuous	<a href="#">Link</a>
	NExT-QA [172]	Video Understanding, Visual QA	Continuous	<a href="#">Link</a>

models' reasoning capabilities. AsDIV [116] and SVAMP [123] focus on word problems, requiring logical reasoning and an understanding of arithmetic relationships in textual contexts. For higher-level reasoning, MATH [55] includes questions in algebra, calculus, and geometry, while AIME targets advanced problem-solving with Olympiad-style questions designed to challenge agents' reasoning depth.

**QA:** QA datasets are designed to evaluate agents' ability to comprehend, reason, and provide accurate answers to questions, often requiring multi-step inference and contextual understanding. HotpotQA [187] is a large-scale multi-hop QA dataset with over 113,000 samples, requiring agents to extract information from multiple documents and perform logical reasoning. StrategyQA [46] focuses on implicit reasoning, requiring multi-step strategies to answer yes/no questions, while MMLU [54] evaluates knowledge across 57 disciplines, testing models on a wide range of zero-shot and few-shot tasks. TruthfulQA [96] tests agents' ability to avoid generating answers based on human-like misconceptions across 38 categories. TriviaQA [70] provides a challenging reading comprehension platform with over 650,000 question-answer-evidence triples, emphasizing linguistic diversity and cross-sentence reasoning. For domain-specific reasoning, PubMedQA [68]

focuses on biomedical questions requiring logical deduction from research abstracts. Datasets like MuSiQue [151] and 2WikiMultihopQA [56] test multi-hop reasoning, integrating structured and unstructured data sources, while QASPER [27] emphasizes ability of extracting answers from academic research papers. ARC [24], which is divided into a challenge set and an easy set, targets grade-school level scientific reasoning, challenging agents to solve multi-choice questions.

**Code:** Code-grounded datasets evaluate agents’ capabilities in code generation and programming tasks, such as translating natural language descriptions into executable code. SWE-bench [184] focuses on real-world software engineering tasks, testing agents on generating patch code for GitHub issues by leveraging natural language problem descriptions and codebase context. Similarly, HumanEval [17] assesses agents’ ability to synthesize Python code from docstrings, serving as a benchmark for fundamental code generation tasks. For broader programming challenges, LiveCodeBench [118] collects problems from competitive platforms such as LeetCode and AtCoder, emphasizing skills like bug-fixing, code execution, and test output prediction. In addition, BIRD [81] and InterCodeSQL [185] evaluate models’ performance in database-related tasks, including generating SQL queries from natural language and refining them through execution feedback. These datasets highlight the complexities of programming environments and the diverse capabilities required for effective code generation and adaptation.

**Tool-use:** Tool-use datasets evaluate agents’ ability to interact with external tools or APIs to accomplish complex tasks. T-Eval [21] provides a fine-grained evaluation framework, breaking tool usage into distinct steps such as instruction following, planning, reasoning, retrieval, and review. Covering 15 fundamental tools across multiple domains, it offers detailed insights into an agent’s performance at each stage of tool utilization. ToolEval [131], based on ToolBench, is an automated evaluator assessing agents’ tool usage capabilities. It uses two metrics: Pass Rate, measuring task completion success, and win rate, comparing preferences between two solution candidates for a given task. MINT-Bench [160] evaluates agents’ ability to handle multi-turn interactions with tools and respond to natural language feedback, emphasizing performance in dynamic tasks. API-Bank [84] tests agents on API planning, retrieval, and execution, covering over 2,000 APIs across 1,000 domains. This dataset assesses agents’ ability to perform real-world tasks such as web search, computation, and smart home control, providing a comprehensive benchmark for tool-augmented models.

**Web:** Benchmarks for web environments evaluate agents’ capabilities in navigating and interacting with online environments to accomplish complex tasks. WebShop [188] simulates an e-commerce platform with 1.18 million products and over 12,000 textual descriptions, challenging agents to understand instructions, reformulate queries, and strategically explore webpages to complete purchases. WebArena [217] offers a comprehensive benchmark with four fully operational web applications, such as online shopping and collaborative development. Additionally, its lightweight version, WebArena-Lite, provides core functionalities for resource-constrained testing. Mind2Web [28] focuses on generalist web agents, covering 137 real-world websites across 31 domains with 2,350 open-ended tasks and crowdsourced operation sequences. MiniWoB++ [97] includes over 100 webpage interaction environments, facilitating research in reinforcement learning and automated web interface agent tasks.

**Environment Interaction:** Environment Interaction tasks evaluate agents’ abilities of complex reasoning, decision-making, and task execution within dynamic environments. These tasks often involve interacting with physical or virtual worlds, requiring agents to adapt and solve problems based on the environment’s feedback. ScienceWorld [157] challenges agents to perform scientific experiments using text-based instructions, requiring logical reasoning and task execution. In the embodied scenarios, ALFWorld [141] and TDW-MAT [201] assess agents’ abilities in physical and multi-agent environments, where tasks include object manipulation, navigation, and collaboration.



In C-WAH [201] and ALFRED [140], agents perform household tasks by interpreting natural language instructions and coordinating actions.

Besides, game-based interaction tasks assess agents’ strategic decision-making and risk assessment abilities in competitive and cooperative settings. RLCard [200] provides several card game environments, such as Blackjack and Limit Texas Hold’em, where agents must develop optimal strategies in incomplete-information games. OpenSpiel [77] offers a variety of games, including the 5x5 Breakthrough game, which tests agents’ ability to devise strategies and optimize policies in multi-agent environments.

**Multimodal:** In multimodal tasks, agents are tested on their ability to integrate and reason with multiple input types, such as images, text, and videos. VQA-V2 [7] is a benchmark for visual QA, consisting of over 265,000 images and 1.1 million questions, challenging agents to answer questions requiring both visual comprehension and textual reasoning. A-OKVQA [137] and ScienceQA (IMG) [108] evaluate agents’ ability to reason across visual and textual information in complex domains like general knowledge and science.

EgoSchema [114] and Next-QA [172] focus on video understanding, with EgoSchema assessing agents’ ability to reason about human activities in long-form video, and Next-QA challenging agents to infer causal and temporal relationships in short videos. Both datasets highlight the complexities of reasoning with time-dependent and contextual information.

**5.1.2 Multi-task Benchmarks.** Multi-task benchmarks evaluate LLM-based agents across various tasks, testing their ability to generalize and adapt to different domains. These benchmarks typically integrate multiple task types into a unified platform, providing a more comprehensive assessment of agents’ performance. Table 8 summarizes key multi-task benchmarks, highlighting the number of datasets covered, task domain types, test sample sizes, and links to the corresponding resources.

Table 8. Multi-Task Agent Evaluation Benchmarks. We use ①, ②, ③, etc., to represent the covered task domains: ① QA, ② Math, ③ Web, ④ Embodied AI, ⑤ Code, ⑥ Tools, ⑦ Game, ⑧ Writing, ⑨ Role-playing, ⑩ Medical Diagnosis.

Benchmark	# Datasets	# Task Types	Test Size	Covered Domains	Data Link
AgentBench [104]	8	3	1,091	③ ⑤ ⑦	<a href="#">Link</a>
AgentEval [170]	14	5	1,160	③ ④ ⑤ ⑥ ⑦	<a href="#">Link</a>
Just-Eval [95]	9	7	1,000	① ② ⑤ ⑧ ⑨	<a href="#">Link</a>
StreamBench [165]	7	5	9,702	① ⑤ ⑥ ⑩	<a href="#">Link</a>
AgentBoard [110]	9	4	1,013	③ ④ ⑥ ⑦	<a href="#">Link</a>

- **AgentBench** evaluates agents across eight environments, including operating systems, databases, and web browsing. It focuses on three task types—programming, games, and web navigation—assessing reasoning, decision-making, and tool usage to provide insights into real-world performance.
- **AgentEval** is a benchmark suite derived from AgentGym [170], including 14 environments, such as web navigation, text games, embodied control, and programming. This dataset features 1,160 diverse instructions collected from collected tasks in different environments, using the unified ReAct format to combine reasoning and action sequences for comprehensive evaluation.
- **Just-Eval** aggregates samples from multiple datasets like AlpacaEval [87], MT-Bench [214], and LIMA [215], evaluating tasks such as text generation and reasoning. A standardized platform is provided for multi-task learning, evaluating agents’ ability to handle diverse inputs and tasks.

- **StreamBench** evaluates agents in continuous improvement settings, testing their ability to adapt and improve over time with tasks like Text-to-SQL generation, Python programming, and medical diagnostics.
- **AgentBoard** provides a multi-round evaluation framework for decision-making and long-term task execution. This framework involves nine task categories, including web, tool use, games, and embodied AI, within 1,013 exemplary environments.

## 5.2 Datasets for Agent Tuning

Fine-tuning datasets are essential for optimizing LLM-based agents, as they provide high-quality trajectories that enable agents to acquire task-specific capabilities and improve performance in targeted scenarios. Although many datasets for evaluation purposes can also be leveraged for fine-tuning, this section focuses on datasets specially designed for agent tuning. These datasets are meticulously constructed to improve agents’ abilities in diverse tasks and environments. In Table 9, we summarize commonly used agent tuning datasets, detailing the number of tasks, covered domains, filtered trajectory counts, and data sources.

Table 9. Datasets for Agent Tuning. We use ①, ②, ③, etc., to represent the task domains as follows: ① QA (Commonsense Reasoning and Open-domain QA), ② Math, ③ Web, ④ Embodied AI, ⑤ Code, ⑥ Tools, ⑦ Game, ⑧ Fact Verification, ⑨ Dialogue, ⑩ General Instructions.

Dataset	# Tasks	# Filtered Trajectories	Covered Domains	Data Link
AgentInstruct [199]	6	1,866	③ ④ ⑤	<a href="#">Link</a>
AgentBank [143]	16	51,287	① ② ③ ④ ⑤	<a href="#">Link</a>
Agent-FLAN [22]	7	24,703	③ ⑤ ⑥	<a href="#">Link</a>
AgentOhana [202]	10	42,600	③ ⑤ ⑥	<a href="#">Link</a>
FireAct [14]	3	1,344	①	<a href="#">Link</a>
ToRA-CORPUS [49]	2	15,538	②	<a href="#">Link</a>
AgentTraj/AgentTraj-L [170]	14	6,130/14,485	③ ④ ⑤ ⑥ ⑦	<a href="#">Link</a>
SMART-Trajectory [197]	17	142,507	① ⑧ ⑨ ⑩	<a href="#">Link</a>

- **AgentInstruct**: A dataset includes over 35,000 instructions derived from six tasks, including ALFWorld [141], WebShop [188], and Mind2Web [28], to cover embodied AI, database operations, and web-based tasks. Trajectories are generated through self-instruction and task derivation, filtered for quality to ensure effective fine-tuning.
- **AgentBank**: Built from 16 public datasets, AgentBank spans five skill dimensions: reasoning, math, programming, web interaction, and embodied intelligence. It provides 51,287 interaction trajectories, supporting both continuous and discrete action spaces.
- **Agent-FLAN**: A restructured dataset for agent tuning, sourced from AgentInstruct [199], ToolBench [131], and ShareGPT [155] datasets. The training corpus is decomposed and redesigned to separate format adherence and general reasoning tasks, ensuring alignment with pre-training domains while preventing overfitting to specific formats.
- **AgentOhana**: Designed to unify trajectories from heterogeneous sources like Mind2Web [28], ToolBench [131], and ToolAlpaca [148], AgentOhana adopts a standardized multi-turn data structure to enable consistent multi-task fine-tuning across diverse agent scenarios.
- **FireAct**: Constructed under the ReAct paradigm, FireAct integrates trajectories generated across three distinct tasks, leveraging multiple prompting methods, including ReAct, CoT, and Reflexion. The dataset captures multi-turn thought-action-observation sequences generated by GPT-4, structured into the ReAct format to improve reasoning, adaptability, and multi-step problem-solving.

- **ToRA-CORPUS**: Specifically crafted for mathematical reasoning, ToRA-CORPUS consists of 16,000 trajectories derived from MATH [55] and GSM8K [26] datasets. Using GPT-4 with greedy decoding and nucleus sampling, the dataset ensures trajectory correctness and robust problem-solving capability.
- **AgentTraj/AgentTraj-L**: Derived from the AgentGym [170], these datasets contain 6,130 and 14,485 filtered trajectories across 14 environments, respectively. AgentTraj provides a foundation for training a generally capable agent, while AgentTraj-L extends the dataset with a larger trajectory set collected through the same pipeline, offering an upper bound for performance using behavioral cloning.
- **SMART-Trajectory**: Focused on long and short trajectories for knowledge-intensive tasks, SAMRT [197] incorporates data from sources such as open-domain QA, commonsense reasoning, and dialogue. It also supports multi-agent frameworks and trajectory learning techniques.

## 6 Application

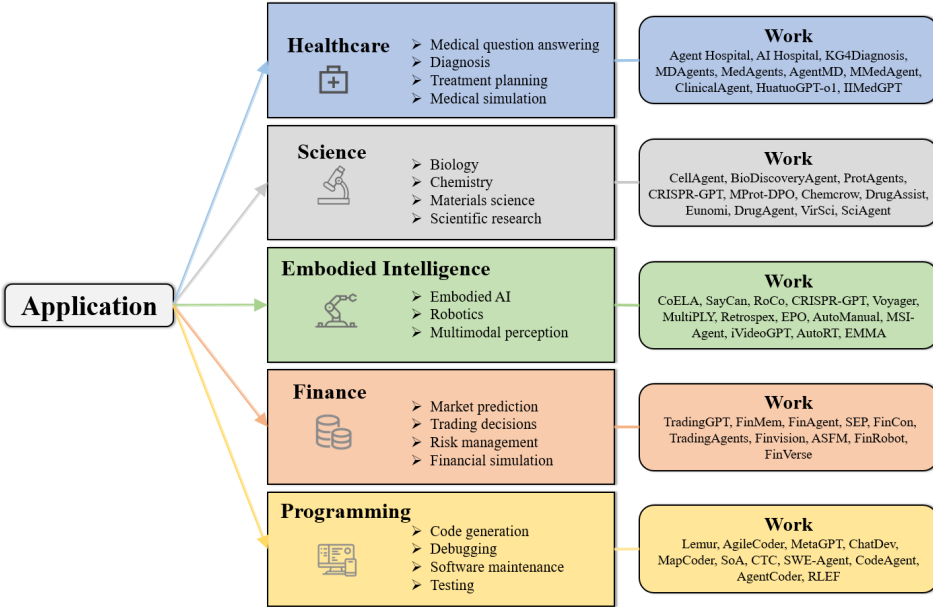


Fig. 4. The applications of LLM-based Agents.

LLM-based agents have been applied across diverse domains, demonstrating their potential to address complex tasks and enhance productivity. This section provides an overview of their roles in healthcare, science, embodied intelligence, finance, and programming, showcasing some representative applications across these fields.

### 6.1 Healthcare

The application of LLM-based agents in healthcare encompasses medical question answering, diagnosis, treatment planning, and medical simulation. Early, LLMs such as Med-PaLM [142], DoctorGLM [176] and BianQue [20], improve conversational accuracy and diagnostic support by fine-tuning medical datasets. DISC-MedLLM [9] further construct medical dataset by using medical

knowledge graphs, real-world dialogue reconstructions, and preference-guided annotations for fine-tuning. Simulation frameworks such as Agent Hospital [82] and AI Hospital [37] introduce dynamic environments, enabling agents to engage in realistic healthcare scenarios involving patients and doctors. KG4Diagnosis [220] and MDAgents [74] facilitate collaborative decision-making by simulating interactions between general practitioners and specialists in multi-agent systems. MedAgents [149] leverages LLM-based agents in a role-playing setting that participate in a collaborative multi-round discussion. Additionally, tool-based approaches like AgentMD [69] and MMedAgent [79] focus on clinical decision support through automated tool selection and multi-modal integration. ClinicalAgent [196] ensembles reasoning technology, multi-agent architectures and external knowledge to boost agent performance in clinical contexts. Recently, HuatuoGPT-o1 [15] and IIMedGPT [209] combine fine-tuning with RL techniques such as PPO and DPO to enhance agents' capabilities in medical reasoning and clinical decision-making. These approaches enable better accuracy in complex tasks, such as generating verified medical responses and optimizing dynamic treatment plans.

## 6.2 Science

AI for Science refers to the application of AI techniques to advance scientific research and discovery. The advent of LLMs has significantly augmented these capabilities, enabling more sophisticated forms of scientific reasoning and automation. LLM-based agents are increasingly utilized across various scientific disciplines, including biology, chemistry, and materials science, to automate data analysis, support hypothesis testing, and optimize research workflows, ultimately accelerating the pace of scientific discovery.

In biological research, LLM-based agents are widely employed for data analysis and experimental design. CellAgent [173] and BioDiscoveryAgent [134] leverage multi-agent collaboration and self-optimization to automate and enhance single-cell data analysis and gene disruption experiments, respectively. Similarly, ProtAgents [47] and CRISPR-GPT [61] utilize LLMs for protein design and gene editing, integrating model inference with external tools to automate complex tasks and expedite biomedical discoveries. Moreover, MProt-DPO [31] employs DPO for protein design, integrating experimental data and simulations to optimize protein landscapes with greater efficacy. In chemistry and materials science, Chemcrow [109] and DrugAssist [191] integrate expert-designed systems to automate organic synthesis, molecule optimization, and chemical problem-solving. These agents enhance LLM capabilities by autonomously planning chemical syntheses and discovering new molecules. Eunomia [6] autonomously extracts and structures datasets from scientific literature, while DrugAgent [102] employs a multi-agent framework for machine learning in drug discovery, optimizing various stages of the research process. In broader scientific research, VirSci [145] and SciAgent [112] demonstrate significant progress in agent-driven discovery. These agents integrate diverse tools to enhance knowledge exchange, improve problem-solving, and accelerate research outcomes.

## 6.3 Embodied Intelligence

Recent advancements in LLM-based agents have significantly enhanced embodied intelligence and robotic systems, enabling autonomously interaction with the physical world. Multimodal inputs such as visual, auditory, and tactile data are integrated in them to perform complex tasks and optimize agent behavior. These agents are usually deployed in physical and text-based environments from household chores to multi-robot coordination, where LLMs guide robots navigating dynamic tasks.

CoELA [201], SayCan [11], and RoCo [113] incorporate agent modules like perception, memory, communication, and planning to perform complex tasks. SayCan combines LLM-driven knowledge

with RL to guide robotic assistants in real-world tasks like kitchen chores, while RoCo supports multi-robot arm collaboration for coordinated actions. In text-based settings, Voyager [156] utilizes long-term memory and reasoning to acquire and refine skills in Minecraft. MultiPLY [58] and Retrospec [171] employ multi-modal sensory inputs and refine action strategies, showcasing how LLM-based agents can bridge perception and planning. Similarly, frameworks like EPO [212], AutoManual [16], and MSI-Agent [41] adopt comparable approaches by breaking down complex tasks into manageable sub-goals, using environmental feedback, and dynamically adapting strategies. Furthermore, iVideoGPT [166], AutoRT [2], and EMMA [186] advance cross-modal capabilities by integrating visual, action, and textual data for task planning and execution. Specifically, iVideoGPT supports interactive video-based reasoning, AutoRT automates multi-robot coordination for real-world data collection, and EMMA transfers LLM expertise across modalities, enabling quick adaptation to new visual tasks.

#### 6.4 Finance

LLM-based agents are widely applied in the financial sector for tasks such as market prediction, trading decision, risk management, and financial simulation [32]. Through advanced reasoning capabilities and external tools, these agents improve decision efficiency, automate repetitive tasks, and adapt to dynamic market environments. TradingGPT [89], FinMem [194], and FinAgent [208] utilize hierarchical memory and reflection mechanisms to analyze historical trading data, facilitating adaptive decisions to market fluctuations. SEP [75] exemplifies RL in trading by incorporating backtesting as an efficient feedback mechanism, where historical market predictions and outcomes serve as rewards to refine agent performance. Multi-agent systems such as FinCon [195], TradingAgents [174], and Finvision [38] simulate real-world team dynamics by integrating structured reports with natural language debates, enabling collaboration among diverse roles like analysts and risk managers. ASFM [43] integrates cross-modal reasoning and simulation frameworks, allowing agents to process multimodal data, evaluate strategies in realistic environments, and adapt to complex financial scenarios. FinRobot [182] proposes a multi-agent AI platform that leverages financial COT, dynamic model configuration, and multi-source LLM integration to enhance financial analysis. Additionally, FinVerse [4] designs an agent system with an embedded code interpreter and employs SFT to train LLMs for versatile financial analysis.

#### 6.5 Programming

The integration of LLM-based agents into programming workflows has grown significantly, assisting with code generation, debugging, and software maintenance [99, 161]. Optimization techniques and external tool usage further improve their reasoning, planning, and tool-use capabilities, allowing better adaptation to complex software development tasks.

Lemur [180] and AgileCoder [119] balance ability of natural language reasoning and coding execution, making LLM-based agents more versatile in dynamic software development environments. MetaGPT [57], ChatDev [128], and MapCoder [66] employ multi-agent collaboration frameworks to enhance software development, structuring workflows where specialized agents handle tasks such as planning, coding, debugging, and testing. Similarly, SoA [65] and CTC [34] introduce scalable collaboration, dynamically adjusting agent roles, and facilitating cross-team communication to optimize large-scale code generation and refinement. Beyond collaboration, LLM-based agents are also optimized for direct interaction with software artifacts. SWE-Agent [183] and CodeAgent [203] enhance repository-wide code navigation, retrieval, and modification by integrating programming tools for automated testing and execution, which allow agents to operate beyond code generation, supporting repository-level understanding and iterative refinement. To further improve adaptability, AgentCoder [60] and RLEF [45] incorporate iterative feedback and execution-driven

learning to refine agent outputs through continuous testing and adjustment, mirroring human-like problem-solving in programming tasks.

## 7 Challenges & Future Directions

We outline the key challenges faced in current state of LLM-based agent optimization and explore potential future research directions.

### 7.1 Robustness to Data Bias

Due to the reliance of LLM-based agents on the quality and consistency of their training data, there is a key challenge that the distribution mismatch between general pre-training data and agent-specific tuning data. Moreover, some datasets often include interactions generated by multi-agent systems or refined by LLMs, which can introduce cognitive biases from the LLMs, potentially amplifying or creating new biases based on prior model behavior. Another issue is bias in automated data filtering and evaluation. LLM-based agents increasingly rely on model-driven filtering and self-assessment to determine which data to retain for training. This can lead to difficulty bias, where the model disproportionately selects tasks it can already solve while discarding more complex examples. Such an automated selection process often classifies tasks according to their solvability, thereby creating a disparity between the training data and real-world deployment scenarios, with the training data typically being less challenging than real-world tasks.

Future research should focus on developing robust techniques to construct datasets [101] and mitigate data bias, such as adversarial training methods, bias testing and exploration of LLM knowledge boundary. Additionally, aligning training data distributions with real-world applications is crucial to reduce bias. Incorporating human-in-the-loop and multi-agent feedback mechanisms could further improve data diversity and quality, enabling better performance in dynamic, complex environments.

### 7.2 Algorithm Adaptability and Efficiency

A key problem in LLM-based agent optimization is balancing algorithmic efficiency with task-specific adaptability. Current methods, such as RL and fine-tuning, often face challenges such as sparse rewards and large action spaces, leading to suboptimal solutions or overfitting to specific data distributions. Algorithms like PPO are effective, but computationally expensive and hard to scale. DPO simplifies the process and is suited for single-step optimization, but most agent tasks require multi-step interactions, exposing a gap in current methods. Moreover, utilizing agent data for LLM fine-tuning may result in overfitting or performance limitations, as it is often constrained by the inherent capabilities of the base model.

Exploring hybrid approaches that combine RL and fine-tuning offers a potential pathway to improving both adaptability and efficiency. Additionally, Leveraging meta-learning and self-supervised learning to develop more generalizable optimization methods could improve agent performance across diverse tasks. Optimizing reward design and integrating efficiency-focused algorithms will also be crucial for reducing computational costs while maintaining high adaptability.

### 7.3 Cross-Domain Adaptation

Cross-domain generalization is essential for the success of LLM-based agents in real-world applications, as agents must adapt to diverse tasks and environments. Current approaches typically focus on improving performance within specific environments, often by using agent trajectory data to fine-tune capabilities for particular domains. While these methods enhance agent performance on targeted tasks, they fail when applied to new unseen domains. This limitation arises from the

agent’s inability to generalize learned knowledge to varying contexts, particularly when there is a significant mismatch between training and real-world data distributions.

To address this limitation, future work should focus on developing more generalizable methods to improve the transferability of LLMs across domains. Techniques like distribution alignment and domain adaptation can enable agents to transfer knowledge more effectively between tasks, even in dynamic environments. Additionally, exploring multi-task learning and domain-invariant representations will enhance an agent’s ability to generalize across diverse scenarios, ensuring robust performance in real-world applications.

#### 7.4 Standardized Evaluation Metrics

Currently, the evaluation of LLM-based agents lacks standardized metrics, making it difficult to compare performance across diverse tasks and domains. Different agents operate in varied environments, such as mathematical reasoning, web navigation, and embodied AI, each relying on distinct evaluation criteria like accuracy, success rate, or reward. This variability makes it difficult to compare optimization methods fairly and assess their generalizability. Multi-agent systems introduce further complexity, as collaborative decision-making lacks unified benchmarks, which hinders comprehensive performance evaluation. Additionally, existing metrics primarily assess task completion rather than extent of optimization, making it challenging to quantify the effectiveness of different enhancement techniques.

Establishing standardized evaluation frameworks provides a promising direction for facilitating fair comparisons across diverse LLM-based agent tasks. This includes developing unified benchmarks that assess optimization effectiveness beyond task completion, incorporating metrics for adaptability, reasoning complexity, and iterative improvement. Furthermore, creating proper metrics that track stepwise optimization progress and integrating preference-based evaluation signals could provide a more comprehensive assessment.

#### 7.5 Parameter-Driven LLM-based Multi-Agent Optimization

Parameter-driven multi-agent optimization remains an underexplored area in LLM-based agent research. Most existing approaches focus on optimizing parameters for single-agent performance, leaving a gap in methods tailored for parameter-driven multi-agent collaboration. Current multi-agent strategies often rely on frozen LLMs, limiting joint optimization and restricting collaborative potential. Without efficient communication, task allocation, and coordination mechanisms, multi-agent systems risk inefficiencies and misalignment. Additionally, optimizing only a single agent’s parameters can lead to over-specialization, conflicting objectives, and suboptimal cooperation, ultimately hindering overall system performance.

Enhancing multi-agent parameter optimization methods is essential for advancing LLM-based systems [197]. Key areas of focus include communication protocols for effective information exchange, distributed learning for adaptation in shared environments, and coordination strategies to enhance collective efficiency. Joint parameter tuning and multi-agent training should be further explored to align agent objectives and maximize system performance. Moreover, reward-sharing mechanisms and hierarchical decision-making could further enhance cooperation, enabling LLM-based multi-agent systems to operate effectively in dynamic and complex environments.

### 8 Conclusion

This survey provides a comprehensive overview of the optimization methods for LLM-based agents, categorizing them into parameter-driven and parameter-free approaches. We focus on parameter-driven optimization, which includes conventional fine-tuning, RL-based methods and hybrid approaches. We first detail the entire process from data construction to fine-tuning in the

conventional fine-tuning part. Subsequently, we outline RL-based optimization, covering both reward function optimization and preference alignment methods. Then, we introduce parameter-free approaches, including feedback mechanisms, historical experience integration, tool usage, RAG, and multi-agent collaboration. Furthermore, we summarize datasets widely utilized to evaluate or fine-tune agents and introduce the diverse real-world applications of LLM-based agents. Finally, we identify several key challenges that need to be addressed and propose future research directions, with the aim of driving the development of more capable and intelligent LLM-based agents.

## References

- [1] Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. 2024. Large language models for mathematical reasoning: Progresses and challenges. *arXiv preprint arXiv:2402.00157* (2024).
- [2] Michael Ahn, Debidatta Dwibedi, Chelsea Finn, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Karol Hausman, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, et al. 2024. Autort: Embodied foundation models for large scale orchestration of robotic agents. *arXiv preprint arXiv:2401.12963* (2024).
- [3] Renat Aksitov, Sobhan Miryoosefi, Zonglin Li, Daliang Li, Sheila Babayan, Kavya Kopparapu, Zachary Fisher, Ruiqi Guo, Sushant Prakash, Pranesh Srinivasan, et al. 2023. Rest meets react: Self-improvement for multi-step reasoning llm agent. *arXiv preprint arXiv:2312.10003* (2023).
- [4] Siyu An, Qin Li, Junru Lu, Di Yin, and Xing Sun. 2024. FinVerse: An Autonomous Agent System for Versatile Financial Analysis. *arXiv preprint arXiv:2406.06379* (2024).
- [5] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403* (2023).
- [6] Mehrad Ansari and Seyed Mohamad Moosavi. 2024. Agent-based learning of materials datasets from the scientific literature. *Digital Discovery* 3, 12 (2024), 2607–2617.
- [7] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*. 2425–2433.
- [8] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection. In *The Twelfth International Conference on Learning Representations*.
- [9] Zhijie Bao, Wei Chen, Shengze Xiao, Kuang Ren, Jiaao Wu, Cheng Zhong, Jiajie Peng, Xuanjing Huang, and Zhongyu Wei. 2023. Disc-medllm: Bridging general large language models and real-world medical consultation. *arXiv preprint arXiv:2308.14346* (2023).
- [10] Xiaohe Bo, Zeyu Zhang, Quanyu Dai, Xueyang Feng, Lei Wang, Rui Li, Xu Chen, and Ji-Rong Wen. 2025. Reflective Multi-Agent Collaboration based on Large Language Models. *Advances in Neural Information Processing Systems* 37 (2025), 138595–138631.
- [11] Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. 2023. Do as i can, not as i say: Grounding language in robotic affordances. In *Conference on robot learning*. PMLR, 287–318.
- [12] Lorenzo Canese, Gian Carlo Cardarilli, Luca Di Nunzio, Rocco Fazzolari, Daniele Giardino, Marco Re, and Sergio Spanò. 2021. Multi-agent reinforcement learning: A review of challenges and applications. *Applied Sciences* 11, 11 (2021), 4948.
- [13] Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2024. ChatEval: Towards Better LLM-based Evaluators through Multi-Agent Debate. In *The Twelfth International Conference on Learning Representations*.
- [14] Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. 2023. Fireact: Toward language agent fine-tuning. *arXiv preprint arXiv:2310.05915* (2023).
- [15] Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, Wanlong Liu, Rongsheng Wang, Jianye Hou, and Benyou Wang. 2024. Huatuogpt-o1, towards medical complex reasoning with llms. *arXiv preprint arXiv:2412.18925* (2024).
- [16] Minghao Chen, Yihang Li, Yanting Yang, Shiyu Yu, Binbin Lin, and Xiaofei He. 2024. AutoManual: Constructing Instruction Manuals by LLM Agents via Interactive Environmental Learning. In *Advances in Neural Information Processing Systems*, Vol. 37. 589–631.
- [17] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).



- [18] Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, et al. 2024. AgentVerse: Facilitating Multi-Agent Collaboration and Exploring Emergent Behaviors. In *ICLR*.
- [19] Weize Chen, Jiarui Yuan, Chen Qian, Cheng Yang, Zhiyuan Liu, and Maosong Sun. 2024. Optima: Optimizing effectiveness and efficiency for llm-based multi-agent system. *arXiv preprint arXiv:2410.08115* (2024).
- [20] Yirong Chen, Zhenyu Wang, Xiaofen Xing, Zhipei Xu, Kai Fang, Junhong Wang, Sihang Li, Jieliang Wu, Qi Liu, Xiangmin Xu, et al. 2023. Bianque: Balancing the questioning and suggestion ability of health llms with multi-turn health conversations polished by chatgpt. *arXiv preprint arXiv:2310.15896* (2023).
- [21] Zehui Chen, Weihua Du, Wenwei Zhang, Kuikun Liu, Jiangning Liu, Miao Zheng, Jingming Zhuo, Songyang Zhang, Dahua Lin, Kai Chen, et al. 2024. T-eval: Evaluating the tool utilization capability of large language models step by step. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 9510–9529.
- [22] Zehui Chen, Kuikun Liu, Qiuchen Wang, Wenwei Zhang, Jiangning Liu, Dahua Lin, Kai Chen, and Feng Zhao. 2024. Agent-FLAN: Designing Data and Methods of Effective Agent Tuning for Large Language Models. *arXiv preprint arXiv:2403.12881* (2024).
- [23] Jihye Choi, Nils Palumbo, Prasad Chalasani, Matthew M Engelhard, Somesh Jha, Anivarya Kumar, and David Page. 2024. MALADE: Orchestration of LLM-powered Agents with Retrieval Augmented Generation for Pharmacovigilance. In *Machine Learning for Healthcare Conference*. PMLR.
- [24] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457* (2018).
- [25] Jesse Clifton and Eric Laber. 2020. Q-learning: Theory and applications. *Annual Review of Statistics and Its Application* 7, 1 (2020), 279–301.
- [26] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168* (2021).
- [27] Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A Smith, and Matt Gardner. 2021. A Dataset of Information-Seeking Questions and Answers Anchored in Research Papers. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 4599–4610.
- [28] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems* 36 (2023), 28091–28114.
- [29] Zhirui Deng, Zhicheng Dou, Yutao Zhu, Ji-Rong Wen, Ruibin Xiong, Mang Wang, and Weipeng Chen. 2024. From Novice to Expert: LLM Agent Policy Optimization via Step-wise Reinforcement Learning. *arXiv preprint arXiv:2411.03817* (2024).
- [30] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems* 36 (2024).
- [31] Gautham Dharuman, Kyle Hippe, Alexander Brace, Sam Foreman, Väinö Hatanpää, Varuni K Sastry, Huihuo Zheng, Logan Ward, Servesh Muralidharan, Archit Vasani, et al. 2024. MProt-DPO: Breaking the ExaFLOPS Barrier for Multimodal Protein Design Workflows with Direct Preference Optimization. In *SC24: International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 1–13.
- [32] Han Ding, Yinheng Li, Junhao Wang, and Hang Chen. 2024. Large language model agent in financial trading: A survey. *arXiv preprint arXiv:2408.06361* (2024).
- [33] Zi-Yi Dou, Cheng-Fu Yang, Xueqing Wu, Kai-Wei Chang, and Nanyun Peng. 2024. Re-rest: Reflection-reinforced self-training for language agents. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 15394–15411.
- [34] Zhuoyun Du, Chen Qian, Wei Liu, Zihao Xie, Yifei Wang, Yufan Dang, Weize Chen, and Cheng Yang. 2024. Multi-agent software development through cross-team collaboration. *arXiv preprint arXiv:2406.08979* (2024).
- [35] Damien Ernst and Arthur Louette. 2024. Introduction to reinforcement learning. *Feuerriegel, S., Hartmann, J., Janiesch, C., and Zschech, P* (2024), 111–126.
- [36] Yue Fan, Xiaojian Ma, Rujie Wu, Yuntao Du, Jiaqi Li, Zhi Gao, and Qing Li. 2024. Videoagent: A memory-augmented multimodal agent for video understanding. In *European Conference on Computer Vision*. Springer, 75–92.
- [37] Zhihao Fan, Lai Wei, Jialong Tang, Wei Chen, Wang Siyuan, Zhongyu Wei, and Fei Huang. 2025. AI Hospital: Benchmarking Large Language Models in a Multi-agent Medical Interaction Simulator. In *Proceedings of the 31st International Conference on Computational Linguistics*. 10183–10213.
- [38] Sorouralsadat Fatemi and Yuheng Hu. 2024. FinVision: A Multi-Agent Framework for Stock Market Prediction. In *Proceedings of the 5th ACM International Conference on AI in Finance*. 582–590.

- [39] Peiyuan Feng, Yichen He, Guanhua Huang, Yuan Lin, Hanchong Zhang, Yuchen Zhang, and Hang Li. 2024. AGILE: A Novel Reinforcement Learning Framework of LLM Agents. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- [40] Xidong Feng, Ziyu Wan, Haotian Fu, Bo Liu, Mengyue Yang, Girish A Koushik, Zhiyuan Hu, Ying Wen, and Jun Wang. 2024. Natural language reinforcement learning. *arXiv preprint arXiv:2411.14251* (2024).
- [41] Dayuan Fu, Biqing Qi, Yihuai Gao, Che Jiang, Guanting Dong, and Bowen Zhou. 2024. MSI-Agent: Incorporating Multi-Scale Insight into Embodied Agents for Superior Planning and Decision-Making. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 643–659.
- [42] Yao Fu, Dong-Ki Kim, Jaekyeom Kim, Sungryull Sohn, Lajanugen Logeswaran, Kyunghoon Bae, and Honglak Lee. 2024. Autoguide: Automated generation and selection of context-aware guidelines for large language model agents. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- [43] Shen Gao, Yuntao Wen, Minghang Zhu, Jianing Wei, Yuhan Cheng, Qunzi Zhang, and Shuo Shang. 2024. Simulating financial market via large language model based agents. *arXiv preprint arXiv:2406.19966* (2024).
- [44] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997* (2023).
- [45] Jonas Gehring, Kunhao Zheng, Jade Copet, Vegard Mella, Taco Cohen, and Gabriel Synnaeve. 2024. Rlef: Grounding code llms in execution feedback with reinforcement learning. *arXiv preprint arXiv:2410.02089* (2024).
- [46] Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics* 9 (2021), 346–361.
- [47] Alireza Ghafarollahi and Markus J Buehler. 2024. ProtAgents: protein discovery via large language model multi-agent collaborations combining physics and machine learning. *Digital Discovery* 3, 7 (2024), 1389–1409.
- [48] Dibya Ghosh, Marlos C Machado, and Nicolas Le Roux. 2020. An operator view of policy gradient methods. *Advances in Neural Information Processing Systems* 33 (2020), 3397–3406.
- [49] Zhibin Gou, Zhihong Shao, Yeyun Gong, yelong shen, Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu Chen. 2024. ToRA: A Tool-Integrated Reasoning Agent for Mathematical Problem Solving. In *The Twelfth International Conference on Learning Representations*.
- [50] Yu Gu, Yiheng Shu, Hao Yu, Xiao Liu, Yuxiao Dong, Jie Tang, Jayanth Srinivasa, Hugo Latapie, and Yu Su. 2024. Middleware for LLMs: Tools Are Instrumental for Language Agents in Complex Environments. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 7646–7663.
- [51] Jian Guan, Wei Wu, zujie wen, Peng Xu, Hongning Wang, and Minlie Huang. 2024. AMOR: A Recipe for Building Adaptable Modular Knowledge Agents Through Process Feedback. In *Advances in Neural Information Processing Systems*, Vol. 37. 126118–126148.
- [52] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).
- [53] Priyanshu Gupta, Shashank Kirtania, Ananya Singha, Sumit Gulwani, Arjun Radhakrishna, Gustavo Soares, and Sherry Shi. 2024. MetaReflection: Learning Instructions for Language Agents using Past Reflections. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 8369–8385.
- [54] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring Massive Multitask Language Understanding. *Proceedings of the International Conference on Learning Representations (ICLR)* (2021).
- [55] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring Mathematical Problem Solving With the MATH Dataset. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks*.
- [56] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing A Multi-hop QA Dataset for Comprehensive Evaluation of Reasoning Steps. In *Proceedings of the 28th International Conference on Computational Linguistics*. 6609–6625.
- [57] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. 2024. MetaGPT: Meta Programming for A Multi-Agent Collaborative Framework. In *ICLR*.
- [58] Yining Hong, Zishuo Zheng, Peihao Chen, Yian Wang, Junyan Li, and Chuang Gan. 2024. Multiply: A multisensory object-centric embodied large language model in 3d world. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 26406–26416.
- [59] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations*.

- [60] Dong Huang, Jie M Zhang, Michael Luck, Qingwen Bu, Yuhao Qing, and Heming Cui. 2023. Agentcoder: Multi-agent-based code generation with iterative testing and optimisation. *arXiv preprint arXiv:2312.13010* (2023).
- [61] Kaixuan Huang, Yuanhao Qu, Henry Cousins, William A Johnson, Di Yin, Mihir Shah, Denny Zhou, Russ Altman, Mengdi Wang, and Le Cong. 2024. Crispr-gpt: An llm agent for automated design of gene-editing experiments. *arXiv preprint arXiv:2404.18021* (2024).
- [62] Xiang Huang, Sitao Cheng, Shanshan Huang, Jiayu Shen, Yong Xu, Chaoyun Zhang, and Yuzhong Qu. 2024. QueryAgent: A Reliable and Efficient Reasoning Framework with Environmental Feedback based Self-Correction. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 5014–5035.
- [63] Xu Huang, Jianxun Lian, Yuxuan Lei, Jing Yao, Defu Lian, and Xing Xie. 2023. Recommender ai agent: Integrating large language models for interactive recommendations. *arXiv preprint arXiv:2308.16505* (2023).
- [64] Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. 2024. Understanding the planning of LLM agents: A survey. *arXiv preprint arXiv:2402.02716* (2024).
- [65] Yoichi Ishibashi and Yoshimasa Nishimura. 2024. Self-organized agents: A llm multi-agent framework toward ultra large-scale code generation and optimization. *arXiv preprint arXiv:2404.02183* (2024).
- [66] Md Ashrafur Islam, Mohammed Eunus Ali, and Md Rizwan Parvez. 2024. MapCoder: Multi-Agent Code Generation for Competitive Problem Solving. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 4912–4944.
- [67] Haolin Jin, Linghan Huang, Haipeng Cai, Jun Yan, Bo Li, and Huaming Chen. 2024. From llms to llm-based agents for software engineering: A survey of current, challenges and future. *arXiv preprint arXiv:2408.02479* (2024).
- [68] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. PubMedQA: A Dataset for Biomedical Research Question Answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2567–2577.
- [69] Qiao Jin, Zhizheng Wang, Yifan Yang, Qingqing Zhu, Donald Wright, Thomas Huang, W John Wilbur, Zhe He, Andrew Taylor, Qingyu Chen, et al. 2024. Agentmd: Empowering language agents for risk prediction with large-scale clinical tool learning. *arXiv preprint arXiv:2402.13225* (2024).
- [70] Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 1601–1611.
- [71] Tomoyuki Kagaya, Thong Jing Yuan, Yuxuan Lou, Jayashree Karlekar, Sugiri Pranata, Akira Kinose, Koki Oguri, Felix Wick, and Yang You. 2024. RAP: Retrieval-Augmented Planning with Contextual Memory for Multimodal LLM Agents. In *NeurIPS 2024 Workshop on Open-World Agents*.
- [72] Dongkyu Kim, Byoungwook Kim, Donggeon Han, and Matouš Eibich. 2024. AutoRAG: automated framework for optimization of retrieval augmented generation pipeline. *arXiv preprint arXiv:2410.20878* (2024).
- [73] Minsoo Kim, Victor Bursztn, Eunyee Koh, Shunan Guo, and Seung-won Hwang. 2024. Rada: Retrieval-augmented web agent planning with llms. In *Findings of the Association for Computational Linguistics ACL 2024*. 13511–13525.
- [74] Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik Siu Chan, Xuhai Xu, Daniel McDuff, Hyeonhoon Lee, Marzyeh Ghassemi, Cynthia Breazeal, Hae Park, et al. 2025. Mdagents: An adaptive collaboration of llms for medical decision-making. *Advances in Neural Information Processing Systems* 37 (2025), 79410–79452.
- [75] Kelvin JL Koa, Yunshan Ma, Ritchie Ng, and Tat-Seng Chua. 2024. Learning to generate explainable stock predictions using self-reflective large language models. In *Proceedings of the ACM Web Conference 2024*. 4304–4315.
- [76] Jakub Lála, Odhran O’Donoghue, Aleksandar Shtedritski, Sam Cox, Samuel G Rodrigues, and Andrew D White. 2023. Paperqa: Retrieval-augmented generative agent for scientific research. *arXiv preprint arXiv:2312.07559* (2023).
- [77] Marc Lanctot, Edward Lockhart, Jean-Baptiste Lespiau, Vinicius Zambaldi, Satyaki Upadhyay, Julien Pérolat, Sriram Srinivasan, Finbarr Timbers, Karl Tuyls, Shayegan Omidshafiei, et al. 2019. OpenSpiel: A framework for reinforcement learning in games. *arXiv preprint arXiv:1908.09453* (2019).
- [78] Dong Won Lee, Hae Park, Yoon Kim, Cynthia Breazeal, and Louis-Philippe Morency. 2024. Global Reward to Local Rewards: Multimodal-Guided Decomposition for Improving Dialogue Agents. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 15737–15762.
- [79] Binxu Li, Tiankai Yan, Yuanting Pan, Jie Luo, Ruiyang Ji, Jiayuan Ding, Zhe Xu, Shilong Liu, Haoyu Dong, Zihao Lin, et al. 2024. MMedAgent: Learning to Use Medical Tools with Multi-modal Agent. In *Findings of the Association for Computational Linguistics: EMNLP 2024*. 8745–8760.
- [80] Dawei Li, Zhen Tan, Peijia Qian, Yifan Li, Kumar Satvik Chaudhary, Lijie Hu, and Jiayi Shen. 2024. Smoa: Improving multi-agent large language models with sparse mixture-of-agents. *arXiv preprint arXiv:2411.03284* (2024).
- [81] Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. 2024. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls.

- Advances in Neural Information Processing Systems* 36 (2024).
- [82] Junkai Li, Yungheui Lai, Weitao Li, Jingyi Ren, Meng Zhang, Xinhui Kang, Siyu Wang, Peng Li, Ya-Qin Zhang, Weizhi Ma, et al. 2024. Agent hospital: A simulacrum of hospital with evolvable medical agents. *arXiv preprint arXiv:2405.02957* (2024).
  - [83] Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. 2024. From Quantity to Quality: Boosting LLM Performance with Self-Guided Data Selection for Instruction Tuning. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. 7595–7628.
  - [84] Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023. A comprehensive benchmark for tool-augmented LLMs. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. 3102–3116.
  - [85] Renhao Li, Minghuan Tan, Derek Wong, and Min Yang. 2024. CoEvol: Constructing Better Responses for Instruction Finetuning through Multi-Agent Cooperation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 4703–4721.
  - [86] Xinyi Li, Sai Wang, Siqi Zeng, Yu Wu, and Yi Yang. 2024. A survey on LLM-based multi-agent systems: workflow, infrastructure, and challenges. *Vicinagearth* 1, 1 (2024), 9.
  - [87] Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. AlpacaEval: An Automatic Evaluator of Instruction-following Models. [https://github.com/tatsu-lab/alpaca\\_eval](https://github.com/tatsu-lab/alpaca_eval).
  - [88] Yunxuan Li, Yibing Du, Jiageng Zhang, Le Hou, Peter Grabowski, Yeqing Li, and Eugene Ie. 2024. Improving Multi-Agent Debate with Sparse Communication Topology. In *Findings of the Association for Computational Linguistics: EMNLP 2024*. 7281–7294.
  - [89] Yang Li, Yangyang Yu, Haohang Li, Zhi Chen, and Khaldoun Khashanah. 2023. Tradinggpt: Multi-agent system with layered memory and distinct characters for enhanced financial trading performance. *arXiv preprint arXiv:2309.03736* (2023).
  - [90] Zaijing Li, Yuquan Xie, Rui Shao, Gongwei Chen, Dongmei Jiang, and Liqiang Nie. 2024. Optimus-1: Hybrid Multimodal Memory Empowered Agents Excel in Long-Horizon Tasks. In *Advances in Neural Information Processing Systems*, Vol. 37. 49881–49913.
  - [91] Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2024. Encouraging Divergent Thinking in Large Language Models through Multi-Agent Debate. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 17889–17904.
  - [92] Xuechen Liang, Meiling Tao, Yinghui Xia, Tianyu Shi, Jun Wang, and JingSong Yang. 2024. Cmat: A multi-agent collaboration tuning framework for enhancing small language models. *arXiv preprint arXiv:2404.01663* (2024).
  - [93] Xuechen Liang, Meiling Tao, Yinghui Xia, Tianyu Shi, Jun Wang, and JingSong Yang. 2024. Self-evolving Agents with reflective and memory-augmented abilities. *arXiv preprint arXiv:2409.00872* (2024).
  - [94] Bill Yuchen Lin, Yicheng Fu, Karina Yang, Faeze Brahman, Shiyu Huang, Chandra Bhagavatula, Prithviraj Ammanabrolu, Yejin Choi, and Xiang Ren. 2024. Swiftsage: A generative agent with fast and slow thinking for complex interactive tasks. *Advances in Neural Information Processing Systems* 36 (2024).
  - [95] Bill Yuchen Lin, Abhilasha Ravichander, Ximing Lu, Nouha Dziri, Melanie Sclar, Khyathi Chandu, Chandra Bhagavatula, and Yejin Choi. 2024. The Unlocking Spell on Base LLMs: Rethinking Alignment via In-Context Learning. In *The Twelfth International Conference on Learning Representations*.
  - [96] Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring How Models Mimic Human Falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 3214–3252.
  - [97] Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. 2018. Reinforcement learning on web interfaces using workflow-guided exploration. *arXiv preprint arXiv:1802.08802* (2018).
  - [98] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. *Advances in neural information processing systems* 36 (2024).
  - [99] 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).
  - [100] Jie Liu, Pan Zhou, Yingjun Du, Ah-Hwee Tan, Cees GM Snoek, Jan-Jakob Sonke, and Efstratios Gavves. 2024. CaPo: Cooperative Plan Optimization for Efficient Embodied Multi-Agent Cooperation. *arXiv preprint arXiv:2411.04679* (2024).
  - [101] Sannyuya Liu, Jintian Feng, Zongkai Yang, Yawei Luo, Qian Wan, Xiaoxuan Shen, and Jianwen Sun. 2024. COMET: “cone of experience” enhanced large multimodal model for mathematical problem generation. *Science China Information Sciences* 67, 12 (2024), 1–2.

- [102] Sizhe Liu, Yizhou Lu, Siyu Chen, Xiyang Hu, Jieyu Zhao, Tianfan Fu, and Yue Zhao. 2024. Drugagent: Automating ai-aided drug discovery programming through llm multi-agent collaboration. *arXiv preprint arXiv:2411.15692* (2024).
- [103] Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022. P-Tuning: Prompt Tuning Can Be Comparable to Fine-tuning Across Scales and Tasks. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. 61–68.
- [104] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. 2024. AgentBench: Evaluating LLMs as Agents. In *The Twelfth International Conference on Learning Representations*.
- [105] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [106] Yun-ting Liu, Jia-ming Yang, Liang Chen, Ting Guo, and Yu Jiang. 2020. Overview of reinforcement learning based on value and policy. In *2020 Chinese Control And Decision Conference (CCDC)*. IEEE, 598–603.
- [107] Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. 2024. A dynamic LLM-powered agent network for task-oriented agent collaboration. In *First Conference on Language Modeling*.
- [108] Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems* 35 (2022), 2507–2521.
- [109] Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. 2024. Augmenting large language models with chemistry tools. *Nature Machine Intelligence* 6, 5 (2024), 525–535.
- [110] Chang Ma, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng Kong, and Junxian He. 2024. AgentBoard: An Analytical Evaluation Board of Multi-turn LLM Agents. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- [111] Hao Ma, Tianyi Hu, Zhiqiang Pu, Boyin Liu, Xiaolin Ai, Yanyan Liang, and Min Chen. 2024. Coevolving with the Other You: Fine-Tuning LLM with Sequential Cooperative Multi-Agent Reinforcement Learning. *arXiv preprint arXiv:2410.06101* (2024).
- [112] Yubo Ma, Zhibin Gou, Junheng Hao, Ruochen Xu, Shuohang Wang, Liangming Pan, Yujiu Yang, Yixin Cao, and Aixun Sun. 2024. SciAgent: Tool-augmented Language Models for Scientific Reasoning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 15701–15736.
- [113] Zhao Mandi, Shreeya Jain, and Shuran Song. 2024. Roco: Dialectic multi-robot collaboration with large language models. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 286–299.
- [114] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. 2024. Egoschema: A diagnostic benchmark for very long-form video language understanding. *Advances in Neural Information Processing Systems* 36 (2024).
- [115] Yuren Mao, Yuhang Ge, Yijiang Fan, Wenyi Xu, Yu Mi, Zhonghao Hu, and Yunjun Gao. 2025. A survey on lora of large language models. *Frontiers of Computer Science* 19, 7 (2025), 197605.
- [116] Shen-Yun Miao, Chao-Chun Liang, and Keh-Yih Su. 2020. A Diverse Corpus for Evaluating and Developing English Math Word Problem Solvers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 975–984.
- [117] Ofir Nachum, Mohammad Norouzi, Kelvin Xu, and Dale Schuurmans. 2017. Bridging the gap between value and policy based reinforcement learning. *Advances in neural information processing systems* 30 (2017).
- [118] King Han Naman Jain, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974* (2024).
- [119] Minh Huynh Nguyen, Thang Phan Chau, Phong X Nguyen, and Nghi DQ Bui. 2024. Agilecoder: Dynamic collaborative agents for software development based on agile methodology. *arXiv preprint arXiv:2406.11912* (2024).
- [120] OpenAI. 2024. GPT-4 Technical Report. arXiv:2303.08774 [cs.CL] <https://arxiv.org/abs/2303.08774>
- [121] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems* 35 (2022), 27730–27744.
- [122] Venkatesh Balavadhani Parthasarathy, Ahtsham Zafar, Aafaq Khan, and Arsalan Shahid. 2024. The ultimate guide to fine-tuning llms from basics to breakthroughs: An exhaustive review of technologies, research, best practices, applied research challenges and opportunities. *arXiv preprint arXiv:2408.13296* (2024).
- [123] Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP Models really able to Solve Simple Math Word Problems?. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2080–2094.
- [124] Jan Peters and Stefan Schaal. 2006. Policy gradient methods for robotics. In *2006 IEEE/RSJ international conference on intelligent robots and systems*. IEEE, 2219–2225.

- [125] Pranav Putta, Edmund Mills, Naman Garg, Sumeet Motwani, Chelsea Finn, Divyansh Garg, and Rafael Rafailov. 2024. Agent q: Advanced reasoning and learning for autonomous ai agents. *arXiv preprint arXiv:2408.07199* (2024).
- [126] Zehan Qi, Xiao Liu, Iat Long Iong, Hanyu Lai, Xueqiao Sun, Wenyi Zhao, Yu Yang, Xinyue Yang, Jiada Sun, Shuntian Yao, et al. 2024. WebRL: Training LLM Web Agents via Self-Evolving Online Curriculum Reinforcement Learning. *arXiv preprint arXiv:2411.02337* (2024).
- [127] Chen Qian, Yufan Dang, Jiahao Li, Wei Liu, Zihao Xie, YiFei Wang, Weize Chen, Cheng Yang, Xin Cong, Xiaoyin Che, Zhiyuan Liu, and Maosong Sun. 2024. Experiential Co-Learning of Software-Developing Agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 5628–5640.
- [128] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, et al. 2024. ChatDev: Communicative Agents for Software Development. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 15174–15186.
- [129] Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, Zhiyuan Liu, and Maosong Sun. 2024. Scaling large-language-model-based multi-agent collaboration. *arXiv preprint arXiv:2406.07155* (2024).
- [130] Shuofei Qiao, Ningyu Zhang, Runnan Fang, Yujie Luo, Wangchunshu Zhou, Yuchen Jiang, Chengfei Lv, and Huajun Chen. 2024. AutoAct: Automatic Agent Learning from Scratch for QA via Self-Planning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 3003–3021.
- [131] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789* (2023).
- [132] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems* 36 (2023), 53728–53741.
- [133] Matthew Renze and Erhan Guven. 2024. Self-Reflection in LLM Agents: Effects on Problem-Solving Performance. *arXiv preprint arXiv:2405.06682* (2024).
- [134] Yusuf Roohani, Andrew Lee, Qian Huang, Jian Vora, Zachary Steinhart, Kexin Huang, Alexander Marson, Percy Liang, and Jure Leskovec. 2024. Biodiscoveryagent: An ai agent for designing genetic perturbation experiments. *arXiv preprint arXiv:2405.17631* (2024).
- [135] Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu, Tianpeng Bao, Guoqing Du, Shiwei Shi, Hangyu Mao, Ziyue Li, Xingyu Zeng, et al. 2023. TPTU: large language model-based AI agents for task planning and tool usage. *arXiv preprint arXiv:2308.03427* (2023).
- [136] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).
- [137] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022. A-okvqa: A benchmark for visual question answering using world knowledge. In *European conference on computer vision*. Springer, 146–162.
- [138] Wentao Shi, Mengqi Yuan, Junkang Wu, Qifan Wang, and Fuli Feng. 2024. Direct Multi-Turn Preference Optimization for Language Agents. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 2312–2324.
- [139] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2024. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems* 36 (2024).
- [140] Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10740–10749.
- [141] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2021. ALFWorld: Aligning Text and Embodied Environments for Interactive Learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- [142] Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023. Large language models encode clinical knowledge. *Nature* 620, 7972 (2023), 172–180.
- [143] Yifan Song, Weimin Xiong, Xiutian Zhao, Dawei Zhu, Wenhao Wu, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. 2024. AgentBank: Towards Generalized LLM Agents via Fine-Tuning on 50000+ Interaction Trajectories. In *Findings of the Association for Computational Linguistics: EMNLP 2024*. 2124–2141.
- [144] Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. 2024. Trial and Error: Exploration-Based Trajectory Optimization of LLM Agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 7584–7600.

- [145] Haoyang Su, Renqi Chen, Shixiang Tang, Zhenfei Yin, Xinzhe Zheng, Jinzhe Li, Biqing Qi, Qi Wu, Hui Li, Wanli Ouyang, Philip Torr, Bowen Zhou, and Nanqing Dong. 2024. Many Heads Are Better Than One: Improved Scientific Idea Generation by A LLM-Based Multi-Agent System. *arXiv preprint arXiv:2410.09403* (2024).
- [146] Hao Sun, Jiayi Wu, Hengyi Cai, Xiaochi Wei, Yue Feng, Bo Wang, Shuaiqiang Wang, Yan Zhang, and Dawei Yin. 2024. AdaSwitch: Adaptive Switching between Small and Large Agents for Effective Cloud-Local Collaborative Learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 8052–8062.
- [147] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. 1999. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems* 12 (1999).
- [148] Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, Boxi Cao, and Le Sun. 2023. Toolalpaca: Generalized tool learning for language models with 3000 simulated cases. *arXiv preprint arXiv:2306.05301* (2023).
- [149] Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark Gerstein. 2024. MedAgents: Large Language Models as Collaborators for Zero-shot Medical Reasoning. In *Findings of the Association for Computational Linguistics ACL 2024*. 599–621.
- [150] Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei Huang, Dacheng Tao, and Jingren Zhou. 2024. A survey on self-evolution of large language models. *arXiv preprint arXiv:2404.14387* (2024).
- [151] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. MuSiQue: Multihop Questions via Single-hop Question Composition. *Transactions of the Association for Computational Linguistics* 10 (2022), 539–554.
- [152] Luong Trung, Xinbo Zhang, Zhanming Jie, Peng Sun, Xiaoran Jin, and Hang Li. 2024. Reft: Reasoning with reinforced fine-tuning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 7601–7614.
- [153] Dennis Ulmer, Elman Mansimov, Kaixiang Lin, Justin Sun, Xibin Gao, and Yi Zhang. 2024. Bootstrapping llm-based task-oriented dialogue agents via self-talk. *arXiv preprint arXiv:2401.05033* (2024).
- [154] Boshi Wang, Hao Fang, Jason Eisner, Benjamin Van Durme, and Yu Su. 2024. LLMs in the Imaginarium: Tool Learning through Simulated Trial and Error. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 10583–10604.
- [155] Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. 2023. Openchat: Advancing open-source language models with mixed-quality data. *arXiv preprint arXiv:2309.11235* (2023).
- [156] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291* (2023).
- [157] Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. 2022. ScienceWorld: Is your Agent Smarter than a 5th Grader?. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 11279–11298.
- [158] Renxi Wang, Haonan Li, Xudong Han, Yixuan Zhang, and Timothy Baldwin. 2024. Learning From Failure: Integrating Negative Examples when Fine-tuning Large Language Models as Agents. *arXiv preprint arXiv:2402.11651* (2024).
- [159] Shenzhi Wang, Chang Liu, Zilong Zheng, Siyuan Qi, Shuo Chen, Qisen Yang, Andrew Zhao, Chaoqi Wang, Shiji Song, and Gao Huang. 2024. Boosting llm agents with recursive contemplation for effective deception handling. In *Findings of the Association for Computational Linguistics ACL 2024*. 9909–9953.
- [160] Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2024. MINT: Evaluating LLMs in Multi-turn Interaction with Tools and Language Feedback. In *The Twelfth International Conference on Learning Representations*. 1–35.
- [161] Yanlin Wang, Wanjuan 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. *arXiv preprint arXiv:2409.09030* (2024).
- [162] Zheng Wang, Zhongyang Li, Zeren Jiang, Dandan Tu, and Wei Shi. 2024. Crafting Personalized Agents through Retrieval-Augmented Generation on Editable Memory Graphs. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 4891–4906.
- [163] Christopher JCH Watkins and Peter Dayan. 1992. Q-learning. *Machine learning* 8 (1992), 279–292.
- [164] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [165] Cheng-Kuang Wu, Zhi Rui Tam, Chieh-Yen Lin, Yun-Nung Chen, and Hung-yi Lee. 2024. StreamBench: Towards Benchmarking Continuous Improvement of Language Agents. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- [166] Jialong Wu, Shaofeng Yin, Ningya Feng, Xu He, Dong Li, Jianye Hao, and Mingsheng Long. 2025. ivideopt: Interactive videopts are scalable world models. *Advances in Neural Information Processing Systems* 37 (2025), 68082–68119.
- [167] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation. *arXiv preprint*

- arXiv:2308.08155* (2023).
- [168] Shirley Wu, Shiyu Zhao, Qian Huang, Kexin Huang, Michihiro Yasunaga, Kaidi Cao, Vassilis Ioannidis, Karthik Subbian, Jure Leskovec, and James Y Zou. 2025. AvaTaR: Optimizing LLM Agents for Tool Usage via Contrastive Reasoning. *Advances in Neural Information Processing Systems* 37 (2025), 25981–26010.
  - [169] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. 2025. The rise and potential of large language model based agents: A survey. *Science China Information Sciences* 68, 2 (2025), 121101.
  - [170] Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Dingwen Yang, Chenyang Liao, Xin Guo, Wei He, et al. 2024. AgentGym: Evolving Large Language Model-based Agents across Diverse Environments. *arXiv preprint arXiv:2406.04151* (2024).
  - [171] Yufei Xiang, Yiqun Shen, Yeqin Zhang, and Nguyen Cam-Tu. 2024. Retrospec: Language Agent Meets Offline Reinforcement Learning Critic. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 4650–4666.
  - [172] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 9777–9786.
  - [173] Yihang Xiao, Jinyi Liu, Yan Zheng, Xiaohan Xie, Jianye Hao, Mingzhi Li, Ruitao Wang, Fei Ni, Yuxiao Li, Jintian Luo, et al. 2024. Cellagent: An llm-driven multi-agent framework for automated single-cell data analysis. *BioRxiv* (2024), 2024–05.
  - [174] Yijia Xiao, Edward Sun, Di Luo, and Wei Wang. 2024. TradingAgents: Multi-Agents LLM Financial Trading Framework. *arXiv preprint arXiv:2412.20138* (2024).
  - [175] Ziyang Xiao, Dongxiang Zhang, Yangjun Wu, Lilin Xu, Yuan Jessica Wang, Xiongwei Han, Xiaojin Fu, Tao Zhong, Jia Zeng, Mingli Song, et al. 2023. Chain-of-experts: When llms meet complex operations research problems. In *The twelfth international conference on learning representations*.
  - [176] Honglin Xiong, Sheng Wang, Yitao Zhu, Zihao Zhao, Yuxiao Liu, Linlin Huang, Qian Wang, and Dinggang Shen. 2023. Doctorglm: Fine-tuning your chinese doctor is not a herculean task. *arXiv preprint arXiv:2304.01097* (2023).
  - [177] Weimin Xiong, Yifan Song, Xiutian Zhao, Wenhao Wu, Xun Wang, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. 2024. Watch Every Step! LLM Agent Learning via Iterative Step-level Process Refinement. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 1556–1572.
  - [178] Fangzhi Xu, Qiushi Sun, Kanzhi Cheng, Jun Liu, Yu Qiao, and Zhiyong Wu. 2024. Interactive Evolution: A Neural-Symbolic Self-Training Framework For Large Language Models. *arXiv preprint arXiv:2406.11736* (2024).
  - [179] Tianyang Xu, Shujin Wu, Shizhe Diao, Xiaozhe Liu, Xingyao Wang, Yangyi Chen, and Jing Gao. 2024. SaySelf: Teaching LLMs to Express Confidence with Self-Reflective Rationales. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 5985–5998.
  - [180] Yiheng Xu, SU Hongjin, Chen Xing, Boyu Mi, Qian Liu, Weijia Shi, Binyuan Hui, Fan Zhou, Yitao Liu, Tianbao Xie, et al. 2024. Lemur: Harmonizing Natural Language and Code for Language Agents. In *The Twelfth International Conference on Learning Representations*.
  - [181] Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. 2024. Large Language Models as Optimizers. In *The Twelfth International Conference on Learning Representations, ICLR 2024*.
  - [182] Hongyang Yang, Boyu Zhang, Neng Wang, Cheng Guo, Xiaoli Zhang, Likun Lin, Junlin Wang, Tianyu Zhou, Mao Guan, Runjia Zhang, et al. 2024. FinRobot: an open-source AI agent platform for financial applications using large language models. *arXiv preprint arXiv:2405.14767* (2024).
  - [183] John Yang, Carlos Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2025. Swe-agent: Agent-computer interfaces enable automated software engineering. *Advances in Neural Information Processing Systems* 37 (2025), 50528–50652.
  - [184] John Yang, Carlos E. Jimenez, Alex L. Zhang, Kilian Lieret, Joyce Yang, Xindi Wu, Ori Press, Niklas Muennighoff, Gabriel Synnaeve, Karthik R. Narasimhan, Diyi Yang, Sida I. Wang, and Ofir Press. 2025. SWE-bench Multimodal: Do AI Systems Generalize to Visual Software Domains?. In *The Thirteenth International Conference on Learning Representations*.
  - [185] John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. 2024. Intercode: Standardizing and benchmarking interactive coding with execution feedback. *Advances in Neural Information Processing Systems* 36 (2024).
  - [186] Yijun Yang, Tianyi Zhou, Kanxue Li, Dapeng Tao, Lusong Li, Li Shen, Xiaodong He, Jing Jiang, and Yuhui Shi. 2024. Embodied multi-modal agent trained by an llm from a parallel textworld. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 26275–26285.
  - [187] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. 2369–2380.



- [188] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems* 35 (2022), 20744–20757.
- [189] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. In *International Conference on Learning Representations (ICLR)*.
- [190] Weiran Yao, Shelby Heinecke, Juan Carlos Niebles, Zhiwei Liu, Yihao Feng, Le Xue, Rithesh Murthy, Zeyuan Chen, Jianguo Zhang, Devansh Arpit, et al. 2023. Retroformer: Retrospective large language agents with policy gradient optimization. *arXiv preprint arXiv:2308.02151* (2023).
- [191] Geyan Ye, Xibao Cai, Houtim Lai, Xing Wang, Junhong Huang, Longyue Wang, Wei Liu, and Xiangxiang Zeng. 2025. Drugassist: A large language model for molecule optimization. *Briefings in Bioinformatics* 26, 1 (2025), bbae693.
- [192] Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. 2024. Agent lumos: Unified and modular training for open-source language agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 12380–12403.
- [193] Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2023. Generate rather than retrieve: Large language models are strong context generators. In *International Conference for Learning Representation (ICLR)*.
- [194] Yangyang Yu, Haochang Li, Zhi Chen, Yuechen Jiang, Yang Li, Denghui Zhang, Rong Liu, Jordan W Suchow, and Khaldoun Khashanah. 2024. FinMem: A performance-enhanced LLM trading agent with layered memory and character design. In *Proceedings of the AAAI Symposium Series*, Vol. 3. 595–597.
- [195] Yangyang Yu, Zhiyuan Yao, Haochang Li, Zhiyang Deng, Yuechen Jiang, Yupeng Cao, Zhi Chen, Jordan Suchow, Zhenyu Cui, Rong Liu, et al. 2025. Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making. *Advances in Neural Information Processing Systems* 37 (2025), 137010–137045.
- [196] Ling Yue, Sixue Xing, Jintai Chen, and Tianfan Fu. 2024. Clinicalagent: Clinical trial multi-agent system with large language model-based reasoning. In *Proceedings of the 15th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics*. 1–10.
- [197] Shengbin Yue, Siyuan Wang, Wei Chen, Xuanjing Huang, and Zhongyu Wei. 2024. Synergistic multi-agent framework with trajectory learning for knowledge-intensive tasks. *arXiv preprint arXiv:2407.09893* (2024).
- [198] Kamer Ali Yuksel and Hassan Sawaf. 2024. A Multi-AI Agent System for Autonomous Optimization of Agentic AI Solutions via Iterative Refinement and LLM-Driven Feedback Loops. *arXiv preprint arXiv:2412.17149* (2024).
- [199] Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. 2023. Agenttuning: Enabling generalized agent abilities for llms. *arXiv preprint arXiv:2310.12823* (2023).
- [200] Daochen Zha, Kwei-Herng Lai, Yuanpu Cao, Songyi Huang, Ruzhe Wei, Junyu Guo, and Xia Hu. 2019. Rlcard: A toolkit for reinforcement learning in card games. *arXiv preprint arXiv:1910.04376* (2019).
- [201] Hongxin Zhang, Weihua Du, Jiaming Shan, Qinhong Zhou, Yilun Du, Joshua B Tenenbaum, Tianmin Shu, and Chuang Gan. 2024. Building Cooperative Embodied Agents Modularly with Large Language Models. In *Proceedings of the International Conference on Learning Representations*.
- [202] Jianguo Zhang, Tian Lan, Rithesh Murthy, Zhiwei Liu, Weiran Yao, Ming Zhu, Juntao Tan, Thai Hoang, Zuxin Liu, Liangwei Yang, et al. 2024. Agentohana: Design unified data and training pipeline for effective agent learning. *arXiv preprint arXiv:2402.15506* (2024).
- [203] Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin. 2024. CodeAgent: Enhancing Code Generation with Tool-Integrated Agent Systems for Real-World Repo-level Coding Challenges. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 13643–13658.
- [204] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. 2021. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of reinforcement learning and control* (2021), 321–384.
- [205] Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. 2023. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792* (2023).
- [206] Shaokun Zhang, Jieyu Zhang, Jiale Liu, Linxin Song, Chi Wang, Ranjay Krishna, and Qingyun Wu. 2024. Offline Training of Language Model Agents with Functions as Learnable Weights. *arXiv:2402.11359 [cs.AI]* <https://arxiv.org/abs/2402.11359>
- [207] Wenqi Zhang, Ke Tang, Hai Wu, Mengna Wang, Yongliang Shen, Guiyang Hou, Zeqi Tan, Peng Li, Yueting Zhuang, and Weiming Lu. 2024. Agent-Pro: Learning to Evolve via Policy-Level Reflection and Optimization. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 5348–5375.
- [208] Wentao Zhang, Lingxuan Zhao, Haochong Xia, Shuo Sun, Jiaze Sun, Molei Qin, Xinyi Li, Yuqing Zhao, Yilei Zhao, Xinyu Cai, et al. 2024. A multimodal foundation agent for financial trading: Tool-augmented, diversified, and generalist. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4314–4325.

- [209] Yiming Zhang, Zheng Chang, Wentao Cai, MengXing Ren, Kang Yuan, Yining Sun, and Zenghui Ding. 2025. IIMedGPT: Promoting Large Language Model Capabilities of Medical Tasks by Efficient Human Preference Alignment. *arXiv preprint arXiv:2501.02869* (2025).
- [210] Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. 2024. A survey on the memory mechanism of large language model based agents. *arXiv preprint arXiv:2404.13501* (2024).
- [211] Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. 2024. Expel: Llm agents are experiential learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 19632–19642.
- [212] Qi Zhao, Haotian Fu, Chen Sun, and George Konidaris. 2024. EPO: Hierarchical LLM Agents with Environment Preference Optimization. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 6401–6415.
- [213] ChuanYang Zheng, Haiming Wang, Enze Xie, Zhengying Liu, Jiankai Sun, Huajian Xin, Jianhao Shen, Zhenguang Li, and Yu Li. 2024. Lyra: Orchestrating Dual Correction in Automated Theorem Proving. *Transactions on Machine Learning Research* (2024).
- [214] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems* 36 (2023), 46595–46623.
- [215] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2023. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems* 36 (2023), 55006–55021.
- [216] Qinhao Zhou, Zihan Zhang, Xiang Xiang, Ke Wang, Yuchuan Wu, and Yongbin Li. 2024. Enhancing the General Agent Capabilities of Low-Paramter LLMs through Tuning and Multi-Branch Reasoning. In *Findings of the Association for Computational Linguistics: NAACL 2024*. 2922–2931.
- [217] Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. 2023. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854* (2023).
- [218] Wangchunshu Zhou, Yixin Ou, Shengwei Ding, Long Li, Jialong Wu, Tiannan Wang, Jiamin Chen, Shuai Wang, Xiaohua Xu, Ningyu Zhang, et al. 2024. Symbolic learning enables self-evolving agents. *arXiv preprint arXiv:2406.18532* (2024).
- [219] Junda Zhu, Lingyong Yan, Haibo Shi, Dawei Yin, and Lei Sha. 2024. ATM: Adversarial Tuning Multi-agent System Makes a Robust Retrieval-Augmented Generator. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 10902–10919.
- [220] Kaiwen Zuo, Yirui Jiang, Fan Mo, and Pietro Lio. 2024. KG4Diagnosis: A Hierarchical Multi-Agent LLM Framework with Knowledge Graph Enhancement for Medical Diagnosis. *arXiv preprint arXiv:2412.16833* (2024).