

TPTU: Large Language Model-based AI Agents for Task Planning and Tool Usage

Jingqing Ruan^{†‡}

ruanjingqing@sensetime.com

Yihong Chen^{†‡}

chenyihong@sensetime.com

Bin Zhang^{†‡}

zhangbin11@sensetime.com

Zhiwei Xu^{†‡}

xuzhiwei@sensetime.com

Tianpeng Bao[†]

baotianpeng@sensetime.com

Guoqing Du[†]

duguoqing@sensetime.com

Shiwei Shi[†]

shishiwei@sensetime.com

Hangyu Mao^{†*}

maohangyu@sensetime.com

Ziyue Li⁺

zlibn@connect.ust.hk

Xingyu Zeng

zengxingyu@sensetime.com

Rui Zhao

zhaorui@sensetime.com

SenseTime Research

Abstract

With recent advancements in natural language processing, Large Language Models (LLMs) have emerged as powerful tools for various real-world applications. Despite their powers, the intrinsic generative abilities of LLMs may prove insufficient for handling complex tasks, which necessitate a combination of task planning and the usage of external tools. In this paper, we first propose a structured framework tailored for LLM-based AI Agents and then discuss the crucial capabilities necessary for tackling intricate problems. Within this framework, we design two distinct types of agents (i.e., one-step agent and sequential agent) to execute the inference process. Subsequently, we instantiate the framework using various LLMs and evaluate their Task Planning and Tool Usage (TPTU) abilities on typical tasks. By highlighting key findings and challenges, our goal is to provide a helpful resource for researchers and practitioners to leverage the power of LLMs in their AI applications. Our study emphasizes the substantial potential of these models while also identifying areas that need more investigation and improvement. The code and resources will be available on GitHub.

1 Introduction

Large Language Model (LLM) [1] is a recent breakthrough in natural language processing (NLP) research. These models are trained on massive amounts of text data and can solve a wide range of tasks, even those that were not included in their training dataset, known as “emerging” ability. This

[†]These authors contribute equally to this work.

[‡]External discussion and ideation.

^{*}These authors work as research interns at SenseTime Research.

⁺The corresponding author.

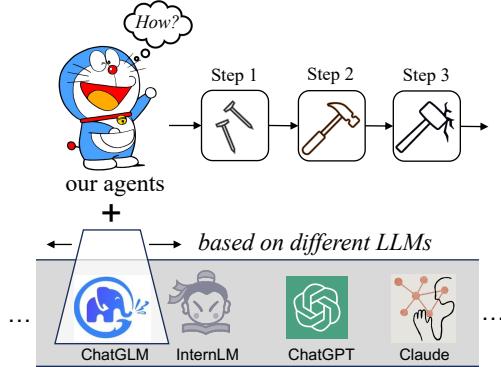


Figure 1: Our LLM-based agents plan tasks and use tools.

ability is especially evident in the tasks of few-shot [2] and zero-shot [3] learning, where LLMs can perform well with minimal or even no fine-tuning to adapt to a new task.

However, the application of LLMs in real-world settings presents unique challenges. On the one hand, LLMs have proved to be incompetent in solving logic problems such as mathematics, and their training data is also out of date (e.g., the knowledge cutoff date for GPT-4 [4] is up to January 2022). Teaching LLMs to use tools such as calculators, calendar, or search engines can help prevent them from hallucinating [5]. On the other hand, despite their impressive problem-solving abilities, the successful integration of these models into complex systems often requires more than just task understanding - it requires the capacity to manipulate various tools and interact effectively with users. This is exemplified in systems like AutoGPT¹, BabyAGI², and ChatGPT-plugins³, which leverage LLMs' capabilities beyond merely generating well-written texts and programs. In these systems, LLMs operate as the central controller, manipulating different tools and interacting with humans, thus taking on the role of Artificial Intelligence Agents (AI Agents). In addition to being central planners, LLMs are often used as intermediaries between macro plans and low-level tool calls or as specific tools. As such, LLMs are seen as a crucial approximation of the linguistic world model in real-world systems.

In this paper, we propose a structured framework for LLM-based AI Agents to evaluate the existing LLMs' planning and tool-using ability and discuss the necessary abilities of such LLM-based AI Agents. Furthermore, we instantiate the framework with different LLMs and evaluate their Task Planning and Tool Usage (TPTU) abilities on several tasks. As shown in Figure 1, we use the Doraemon as an analogy of our LLM-based agents: Doraemon's magic 4D pocket consists of millions of gadgets (the Tool Set), and Doraemon needs to pick the right tools and solve tasks in a right order. Our main contributions are summarized as follows:

1. We propose a structured framework tailored for LLM-based AI Agents to evaluate the TPTU abilities of the existing open-source LLMs.
2. We design two distinct types of agents, namely, one-step agent and sequential agent, to execute the inference process of conducting sub-tasks in a once-for-all or sequential manner, respectively. We provide detailed empirical results and analysis.
3. Our study reveals significant potential in utilizing LLMs for complex tasks. Furthermore, we conclude four following potential weaknesses of LLM-based agents: failing to output in a specific format, struggling to grasp task requirements, over-utilizing one tool, and lack of summary skills. These observations could spark some insights and shed light on the areas that deserve further investigation and improvement.

¹<https://github.com/Significant-Gravitas/Auto-GPT>

²<https://github.com/yoheinakajima/babyagi>

³<https://openai.com/blog/chatgpt-plugins>

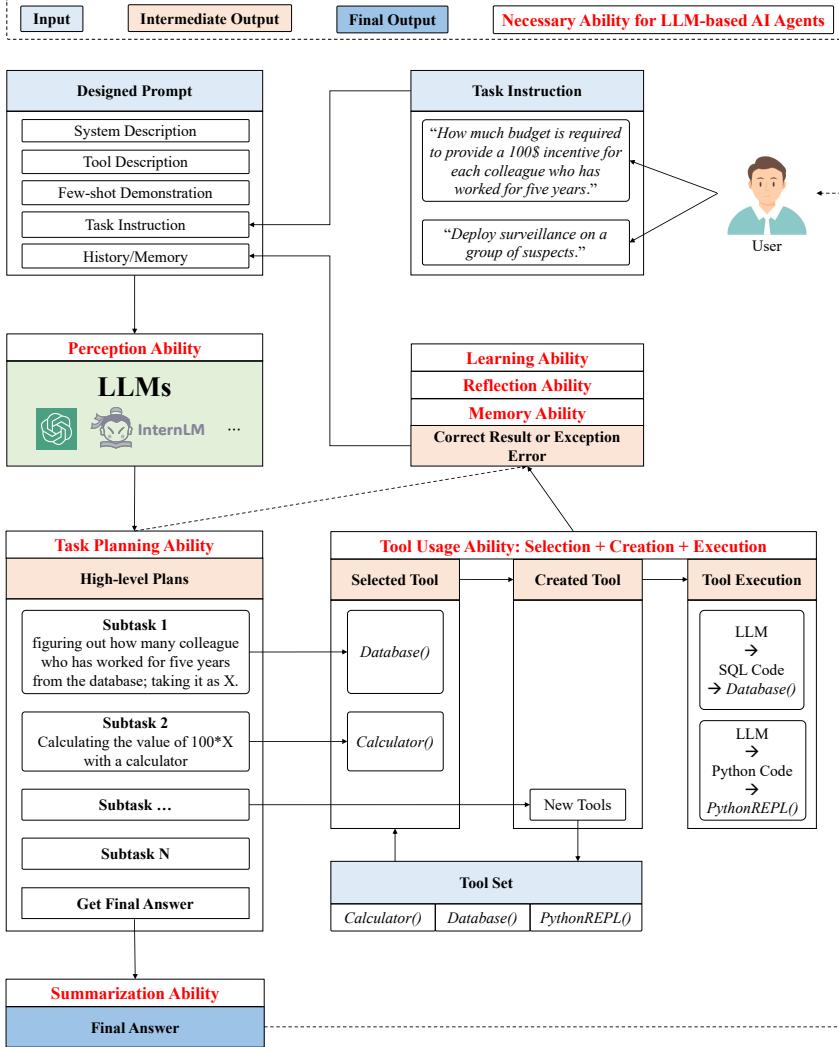


Figure 2: The proposed framework for LLM-based AI Agents.

2 Method

To the best of our knowledge, the study of “Agent”, “Autonomous Agent”, “AI Agent” and “Multi-Agent” has been a central part of AI research for decades [6–11], aimed at understanding and building intelligent and autonomous systems, but there is currently no standardized definition for AI Agents, particularly those that are based on LLMs.

In this paper, the Artificial Intelligence Agent (AI Agent) is defined as a program that employs AI techniques to perform tasks that typically require human-like intelligence. AI Agents can take many forms, from simple chatbots to complex autonomous systems that interact with their environment and make decisions in real-time. They can be trained using a variety of machine learning techniques, including supervised, unsupervised, and reinforcement learning, and can be programmed to perform specific tasks or learn from their experiences in order to improve their performance over time.

2.1 Agent Framework

We are particularly interested in the AI Agent that employs the LLM techniques (i.e., LLM-based AI Agent), due to its high efficiency and flexibility in various tasks and domains. Specifically, we design our AI Agent framework with six components as shown in Figure 2:

- Task Instruction.** This is the explicit input of the agent. In practical systems, the task instruction comes from human users of the systems. For example, in a human resources (HR) system, the user may give a task instruction: How much budget is required to provide a 100\$ incentive for each colleague who has worked for five years? In contrast, in a criminal investigation system, the user may give a task instruction: Deploy surveillance on a group of suspects.
- Designed Prompt.** This is an additional form of input for the agent, derived from tasks that the human users anticipate the AI Agent will complete. Humans can craft specific instructions or demonstrations to steer the LLM-based AI Agents toward generating suitable responses. These guiding inputs could encompass system instructions, tool descriptions, few-shot demonstrations, chat history, or even error output.
- Tool Set.** It is another input for the agent, which refers to the set of external resources, services, or subsystems that the AI Agent can utilize to aid in its tasks. This could include databases for information retrieval [12], APIs for interacting with external systems [5], other AI models specialized for tasks such as image recognition or sentiment analysis [13], or even non-AI tools and resources such as web scraping tools or data visualization libraries [14]. The toolset expands the capabilities of the AI Agent, enabling it to access and process information beyond its internal knowledge, interact with other systems, or perform specialized tasks that it may not be capable of on its own. For example, an AI Agent might use a weather API to fetch current weather information, or a Python interpreter to solve the mathematical question.
- LLM.** This is the core component of the system that interprets the task instructions and prompts, interacts with the toolset, and generates the intermediate outputs and final answers. In this context, we utilize publicly available large language models such as ChatGPT, GPT-4 [4], InterLM [15], and others.
- Intermediate Output.** This represents the output generated by the LLM-based AI Agent after it processes the task instructions and prompts, and interacts with the toolset. There are three typical intermediate outputs: (1) the high-level plans to fulfill the original user instruction, (2) selected and created tools to fulfill each subtask in the plans, and (3) the results or errors produced after tool execution. The output can be reviewed and refined, either by the AI Agent itself or with human oversight, to ensure it is accurate and meets the requirements of the task instruction.
- Final Answer.** This is the output that the AI Agent summarizes and provides to the user after all processing (including task planning, tool usage, and possibly error feedback) has been completed.

2.2 Agent Ability

To apply LLM-based AI Agents to augment or replace human decision-making in real-world applications, the agents typically require the following abilities:

- Perception Ability:** AI Agents must be able to perceive the task instruction from human and system specifications.
- Task Planing Ability:** AI Agents should have the capacity to create a step-by-step plan for complex task composition based on the perceived instruction and specifications. This usually involves the generation of critical subtask sequences, and the ability to adjust the plan dynamically in response to changes in the task or environment.
- Tool Usage Ability:** On the one hand, AI Agents should possess the capacity to **select** a variety of existing tools or resources to execute complex tasks. On the other hand, AI Agents should **create** new tools by converting the task requirements. This ability enables the AI Agent to extend its capabilities beyond LLM itself and the existing tools by leveraging the vast resources available in the digital world. Finally, AI Agents should be able to **execute** the selected or created tools for truly grounding the human request based on the resources and constraints of systems.
- Learning/Reflection/Memory (from Feedback):** AI Agents should be capable of learning from feedback, including correct results and exception errors. They should incorporate

memory, such as logging or chat history, and reflection to adapt their plans or decisions. This allows the agents to improve their performance and efficiency in task execution continuously.

5. **Summarization:** After several rounds of interaction with humans, tools, and systems, AI agents can ultimately complete the original task provided by the users. At this point, AI agents should be able to summarize the interaction history and provide a final answer that is concise and easy to understand for the users.

To endow AI Agents with the aforementioned abilities, some techniques that can be used include chain-of-thought (CoT) and vector databases, as shown in Table 1.

Table 1: A simple illustration of the techniques for endowing the key ability.

Ability	Possible Techniques
Perception	Multi-input Fusion
Task Planing	Zero-shot CoT and Few-shot CoT
Tool Usage (Selection/Creation/Execution)	Text Matching/Code Generation/ Action Grounding
Learning/Reflection/Memory	RLHF/Multi-agent Debate/ Vector Database
Summarization	Attention Mechanism and Natural Language Generation

2.3 Agent Design

Task planning and tool usage represent the cornerstone of LLM's abilities. Others like perception, learning/reflection/memory (from feedback), and summarization are indeed critical, but they primarily serve to enhance and support these two core competencies. Therefore, concentrating on these two key competencies - Task Planning and Tool Usage (TPTU for short) - we have devised two distinct types of AI agents, as depicted in Figure 3:

- The first one, named as the **One-step Agent (TPTU-OA)**, adopts a global perspective to interpret the original problem, effectively breaking it down into a sequence of sub-tasks in a single instance. This strategy fully harnesses the model's comprehensive understanding capabilities to map out the problem-solving steps for all sub-tasks at once. This method underscores the significance of a holistic understanding and planning of the overall task, albeit it might lack flexibility when dealing with individual sub-tasks.
- The second type, referred to as the **Sequential Agent (TPTU-SA)**, emphasizes tackling the current sub-task at hand. Upon successfully resolving the ongoing sub-task, this agent requests the LLMs to provide the succeeding sub-task. This approach enables the model to maintain a clear and concentrated focus throughout the problem-solving journey, tackling issues incrementally. Such a methodology allows for continuous feedback and progress within the confines of addressing a broader problem.

These two distinct agent models represent two disparate problem-solving strategies - the one-step and sequential resolution⁴. In our subsequent experiments, we aim to understand their respective strengths and weaknesses and how they can be best utilized to leverage the capabilities of LLMs in real-world problem-solving scenarios.

3 Evaluation

We instantiate the proposed LLM-based AI Agent framework (TPTU-OA and TPTU-SA) with different LLMs and evaluate their performance on typical tasks.

⁴One can also combine the two strategies to design a hierarchical agent, but this is beyond the scope of this paper.

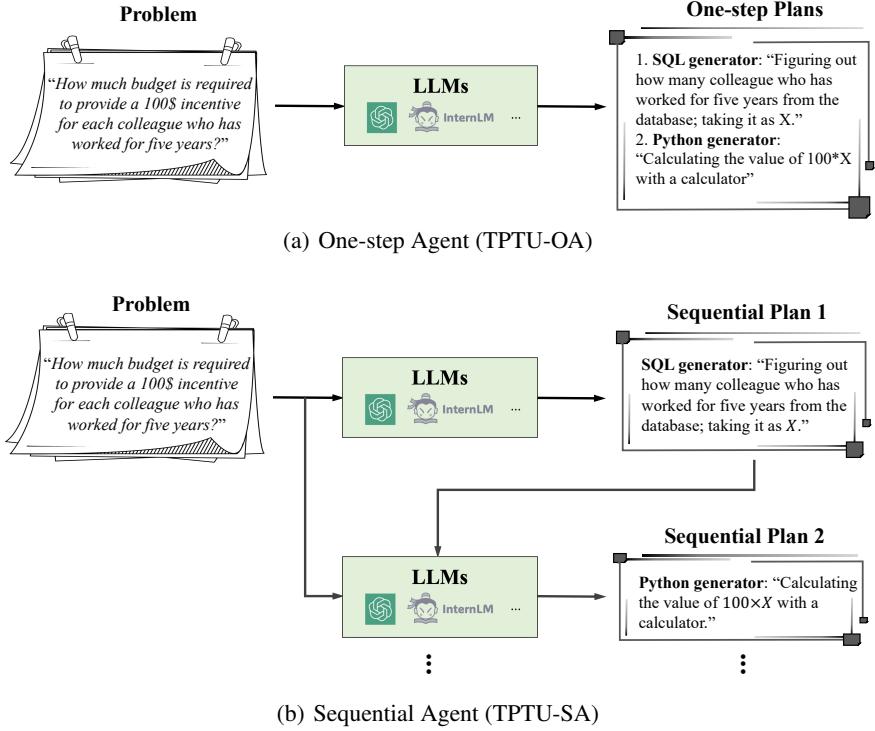


Figure 3: The workflows of the One-step Agent and the Sequential Agent are specifically designed to assess the Task Planning and Tool Usage abilities of LLMs.

3.1 Preparations

Before beginning our evaluation, we first outline the preparations. We will give detailed descriptions of the datasets, available tools, and popular large language models.

3.1.1 Datasets

We first clarify the motivations behind our choice of tools for evaluation. The selection was guided by two primary factors: **the number of tools** to be evaluated and **the specific tools** to be included.

Firstly, regarding the number of tools, it is important to state that our proposed evaluation framework is extensible. It can incorporate any number of tools as pluggable components to be managed by the LLM-based AI agents. Secondly, looking at the current work on tool-augmented LLMs, such as T-Bench [16] and ToolBench [17], we see that only a handful of tools are launched and executed in a single scenario. Meanwhile, API-Bank [18], in a single scenario, typically dispatches only one API tool and awaits its response. APIBench [19] and ToolAlpaca [20] do not even execute a tool response. Hence, for the sake of simplicity and focus in our evaluation, we have decided to primarily assess two tools (which can be called multiple times) within a single scenario.

Secondly, we also need to decide which specific tools should be used for evaluation. Consider a real-world scenario where we pose the question: "How much budget is required to offer a \$100 incentive to each employee who has been with the company for over five years?". To answer this, we first need to retrieve the relevant data from a database, typically using SQL, to find the number of eligible employees. Then, we need to perform a mathematical calculation to estimate the total budget. Such scenarios are quite common in daily life where the formulation and resolution of a question often involve SQL and mathematical tools.

Recognizing the importance of these tools, we have chosen to focus our evaluation on SQL and Python generators, which represent the capabilities of database querying and mathematical computation, respectively. To this end, we have prepared 120 question-answer pairs that vary in complexity. These pairs provide a rigorous assessment of the LLM-based AI agents in understanding, generating, and

utilizing these essential tools. For further information on these queries and their corresponding demonstrations, please refer to Appendix A.

3.1.2 Tools

We have defined a total of 12 available tools for the selection of the LLM-based AI agents for evaluation. They are defined as follows:

- **SQL generator:** Given an input question and a database, create a syntactically correct SQLite query statement.
- **Python generator:** Given an input question and some information, generate a syntactically correct Python code.
- **Weather query tool:** Given a location, output the current real-time weather at that location.
- **Image generator:** Given a text description, generate a related image.
- **Text extractor:** Given a link to an image, extract the corresponding text and its position coordinates.
- **Translator:** Given a piece of text, translate it into other languages.
- **Bing Searcher:** Given a piece of text, conduct a search on the Bing browser and return content.
- **Shell generator:** Given an input question and some information, generate a syntactically correct Shell code.
- **Java generator:** Given an input question and some information, generate a syntactically correct Java code.
- **Wikipedia searcher:** Given a piece of text, conduct a search on Wikipedia and return content.
- **Office software:** Given a text description, automatically generate corresponding long documents or spreadsheets or PPTs.
- **Movie player:** Given a movie name, automatically play the corresponding movie resources.

3.1.3 LLMs

The LLMs evaluated in this paper are listed in Table 2, elaborated as follows:

- **GPT** series developed by OpenAI boasts a powerful language model with a vast number of parameters, enabling it to tackle intricate problems efficiently. This paper aims to evaluate the performance of ChatGPT, which balances the performance with costs (the number of OpenAI API calls).
- **Claude** is committed to maintaining honesty and ensuring user safety, which is developed by Anthropic. With its impressive size, Claude ranks among the largest language models globally and poses a formidable challenge to ChatGPT as a strong competitor.
- **InternLM**, a sophisticated language model developed by Shanghai AI Lab, boasts a multi-round dialogue capability and an impressive ability to comprehend super-long text. This language model is meticulously designed to cater to the nuances of the Chinese language, enabling it to comprehensively understand and effectively process Chinese text. Here, we adopted the version with 120 billion parameters.
- **Ziya** is an expansive and robust pre-training model developed by IDEA, derived from the LLaMa with 13 billion parameters. This comprehensive model exhibits a wide range of capabilities, including translation, programming, and mathematical calculations. Notably, it stands out as a bilingual LLM, highlighting its ability to effectively process and comprehend text in Chinese.
- **ChatGLM**, developed by Tsinghua University, is an open-source dialogue language model that supports bilingual Q&A in Chinese and English, with a particular focus on Chinese optimization. Built on the General Language Model (GLM) architecture and utilizing model quantization technology, the ChatGLM can be easily deployed on consumer-grade graphics cards, enabling local implementation by users.

- **Chinese-Alpaca-Plus** is achieved by extending LLaMA’s existing vocabulary with an additional 20,000 Chinese tokens from Meta AI (formerly known as Facebook AI Research Laboratory). In this version, we use a model with 33 billion parameters. The training text has been expanded to 120GB, and the fine-tuning instruction data has been increased to 4.3M.

Table 2: The LLMs evaluated in this paper.

Organization	Model Name	Model Parameters
OpenAI	ChatGPT[21]	200B
Anthropic	Claude[22]	>52B
Shanghai AI Lab	InternLM	120B
IDEA	Ziya-13B	13B
Tsinghua University	ChatGLM-130B[23]	130B
-	Chinese-Alpaca-Plus-33B[24, 25]	33B

3.2 Evaluation on Task Planning Ability

In this section, to evaluate the planning capabilities of the LLM-based AI agents, we have structured the evaluations as follows.

For TPTU-OA, we begin by examining the agents’ ability to plan the order of tool use. This is followed by an evaluation of the agents’ capacity to not only plan the sequence of tools but also the corresponding subtask descriptions. Subsequently, we conduct a specialized planning evaluation where the agents must generate multiple sequences of key-value pairs of the form {tool: subtask description} in complex problem teardowns. Moreover, we expand the toolset with additional, unrelated tools to further challenge and reassess the planning ability of the LLM-based AI agents.

For TPTU-SA, we follow the regime that the agent should generate multiple sequences of key-value pairs of the form {tool: subtask description} for evaluation.

3.2.1 TPTU-OA: Tool Order Planning

Here, we utilize two kinds of tools for problem-solving: the SQL generator, which retrieves data from databases, and the Python generator, adept at addressing mathematical questions.

To validate the capacity of the LLM-based AI agents to strategically plan for the tool order, we designed the prompt as shown in Figure 8 of Appendix B. This design is motivated by the goal to assess the ability of LLM-based AI agents to understand complex problems, subsequently decomposing them into a sequence of simpler tasks executed by appropriately selected tools. Specifically, we require the LLM-based AI agent to follow our instructions, select tools from our pre-defined tool set with detailed function descriptions, conform to the given format strictly, and understand the demonstrations to learn from them.

Upon feeding these prompts into the LLM-based AI agents under evaluation, we obtained the following accuracy rates for the tool planning, as shown in Table 3.

Table 3: The evaluation results for the planning of tool order generation.

Model	ChatGPT	Claude	Ziya
Accuracy	100%	100%	45%
Model	ChatGLM	Chinese-Alpaca-Plus	InternLM
Accuracy	45%	20%	80%

The results of our experiments indicate that models, notably Ziya and ChatGLM, frequently grapple with the generation of lists in the correct format. For other models, the predominant challenges lie in

generating tools in the correct sequence or in the occasional omission of necessary tools. Nonetheless, the issue of parsing list formats is generally negligible.

These findings suggest that the majority of LLM-based AI agents possess a fundamental capability to analyze the tool needs of a given problem and understand its task requirements. To further explore whether these LLM-based AI agents can effectively break down the original problem into sub-tasks, we proceed to the following section.

3.2.2 TPTU-OA: Tool Order Planning and Subtask Description Generation

Simply planning the order of tool usage is not sufficient to fully address a problem. To truly solve it, we need to provide a guide or instructions for the usage of each tool, that is, a decomposed subtask description. Therefore, we can decompose the original complex problem into two separate sequences. One sequence represents the order in which the tools are utilized, while the other sequence corresponds to the subtask descriptions that each tool in the tool sequence aims to resolve. A problem is only truly solved when both the tool and subtask description sequences have been successfully planned. In order to verify whether LLM-based AI agents truly have the ability to solve complex problems, we designed a new prompt as shown in Figure 9 of Appendix B. The main improvement is to plan the corresponding subtask description for each tool after the tool planning is completed.

Table 4: The evaluation results for the planning of tool order and subtask description generation.

Model	ChatGPT	Claude	Ziya
Accuracy	55%	15%	10%
Model	ChatGLM	Chinese-Alpaca-Plus	InternLM
Accuracy	10%	0%	45%

After feeding the prompt to these LLM-based AI agents, we get results shown in Table 4.

Although the generation of tool sequences and their corresponding subtask descriptions might be an effective way to problem-solving, there is a significant decrease in accuracy for all LLMs as can be seen. We hypothesize that there are a few potential drawbacks to this method:

1. **Difficulty in Error Tracking and Debugging.** Generating the complete tool and subtask sequences may make it more challenging to track and debug errors. If an error arises within the sequence, it might require a total regeneration instead of a simple modification or repair to the erroneous part.
2. **Tool-Subtask Pairing Issue.** If all tool sequences and subtask descriptions are generated independently, there's an inherent risk of misalignment between the tools and their corresponding subtasks. This could potentially lead to an improper pairing, which, in turn, could result in a flawed or ineffective solution that fails to appropriately resolve the given problem.
3. **Lack of Flexibility.** The approach may lack this flexibility when facing complex problems requiring adjustments to the tool or subtask sequence.
4. **Dependency on Global Information.** Generating the entire tool and subtask sequences requires a global understanding and planning of the entire problem. However, in some instances, certain parts of the problem might not be clear at the early stages of problem-solving, which could pose challenges within this framework.

3.2.3 TPTU-OA: The Planning of Tool-Subtask Pair

To mitigate the aforementioned issue, we propose a novel approach to foster flexible problem-solving with the LLM-based AI agent. We prompt the agent to generate multiple sequences, each consisting of a key-value pair in the format of {tool: subtask description} that associates a tool with its respective subtask description. This allows us to simultaneously plan the tool choice and subtask without the risk of improper matching. Moreover, it offers the flexibility to update the planned sequences in real-time based on evolving problem feedback, enhancing adaptability and efficiency when addressing complex tasks.

With this consideration, we have designed a unique prompt that encourages this advanced problem-solving strategy. In the following section, we delve into the specifics of this prompt design in Figure 10 of Appendix B. The key improvement in this prompt is its directive for the LLM-based AI agents to stringently adhere to the predefined dictionary format. To facilitate this, we offer several demonstrations in our desired format, serving as references for the language model to follow.

Table 5: The evaluation results for the planning of Tool-Subtask pair.

Model	ChatGPT	Claude	Ziya
Model	ChatGLM	Chinese-Alpaca-Plus	InternLM
Accuracy	75%	90%	20%
Model	ChatGLM	Chinese-Alpaca-Plus	InternLM
Accuracy	0%	5%	55%

After feeding the prompt to these LLM-based AI agents, we get results shown in Table 5.

Analyzing the results from Tables 4 and 5, we observe a marked improvement of 52.9% when the tool-subtask pairs are generated in a unified format compared to separate generation of tools and subtasks.

This significant performance enhancement can likely be attributed to the close coupling between tools and their associated subtasks in our unified generation strategy. When tools and subtasks are generated separately, there is a potential disconnect or lack of coherence between the two, which could lead to less accurate or efficient solutions. In contrast, by generating tool-subtask pairs together, we ensure that each tool is directly tied to its relevant subtask, leading to a more coordinated and effective problem-solving approach. This might explain the observed increase in overall performance.

3.2.4 TPTU-OA: The Planning of Tool-Subtask Pair with Unrelated Tools

So far, our analysis and evaluation have been primarily focused on the LLM-based AI agents' proficiency in planning with specific tools. However, we are also interested in how it would perform when faced with many irrelevant or similar tools. Therefore, for a more comprehensive assessment, we expanded the prompt in Table 10 to include an additional ten unrelated tools, as illustrated in Figure 11 of Appendix B.

Table 6: The evaluation results for the planning of Tool-Subtask pair with unrelated tools.

Model	ChatGPT	Claude	Ziya
Model	ChatGLM	Chinese-Alpaca-Plus	InternLM
Accuracy	70%	90%	10%
Model	ChatGLM	Chinese-Alpaca-Plus	InternLM
Accuracy	0%	5%	50%

After feeding the prompt to these LLM-based AI agents, we get results shown in Table 6. The results from our expanded evaluation demonstrate that even when presented with irrelevant or similar tools and descriptions, LLM-based AI agents consistently avoid selecting these unrelated tools (i.e., the accuracy has remained unchanged or exhibited only a marginal decrease compared with Table 5). This outcome indicates the effectiveness of our designed prompt, which successfully guides the LLM-based agents to understand the appropriate tool sequence for complex problem decomposition.

This observation reinforces the notion that a well-structured and informative prompt can efficiently guide AI agents to understand the core essence of the problem, thereby enabling them to sift through irrelevant information and focus on key tasks. This successful discrimination against unrelated tools also points towards the models' ability to understand the specific context of a problem and select the appropriate tools, thereby enhancing the overall problem-solving process.

3.2.5 TPTU-SA: The Planning of Tool-Subtask Pair Generation

Upon identifying the drawbacks of first generating a list of tools and then generating corresponding subtask descriptions, we decided to focus subsequent tests on the generation of tool-subtask pairs.

Consequently, in this section, we evaluate the capability of TPTU-SA to generate these tool-subtask pairs.

To achieve the goal of recursively generating tool-subtask pairs, we have designed prompts as illustrated in Figure 12 of Appendix B.

Table 7: The evaluation results for the planning of Tool-Subtask with the sequential agent.

Model	ChatGPT	Claude	Ziya
Accuracy	80%	100%	10%
Model	ChatGLM	Chinese-Alpaca-Plus	InternLM
Accuracy	0%	0%	65%

The evaluation results are shown in Table 7. Compared with results shown in Table 5, TPTU-SA generally performs better than TPTU-OA especially for high-performing LLMs (e.g., ChatGPT, Claude and InternLM). We propose the following potential reasons for this observation:

1. **Sequentiality Mimics Human Problem-Solving:** In real-world scenarios, humans tend to solve complex problems by breaking them down into smaller, manageable subtasks which are often handled sequentially. Sequential agents are designed to mimic this step-by-step approach, which might inherently suit complex problem-solving better.
2. **Richer Contextual Understanding:** Sequential agents are exposed to the outcome of each previous subtask before moving on to the next one. This iterative process could facilitate a richer understanding of the problem context, enabling more accurate task planning and tool usage.
3. **Flexibility in Task Management:** In comparison to one-step agents, sequential agents might have more flexibility in managing tasks. They have the opportunity to correct errors or adjust their strategy after each step, which can lead to improved overall performance.
4. **Improved Learning From History:** The sequential process provides a history of actions and results which can be beneficial in learning. The agent can use this history to make better predictions about what tool to use next or what subtask to tackle, leading to more accurate and efficient problem-solving.

These points of analysis suggest that the structure and operation of sequential agents inherently confer certain advantages in complex problem-solving scenarios, leading to their superior performance.

3.3 Evaluation on Tool Usage Ability

Before evaluating the end-to-end multi-tool usage ability of LLM-based AI agents, we first evaluate the effectiveness of single-tool usage for SQL generation and mathematical code generation.

Subsequently, to assess the end-to-end performance of LLMs across various tools, two types of agents (TPTU-OA and TPTU-SA) were developed and several LLMs were subjected to testing under these agents. The role of the agents is to break down complex questions into simpler sub-questions and plan corresponding tools to solve them, based on the available toolset and corresponding tool descriptions.

3.3.1 The effectiveness of Single Tool Usage

Our aim is to systematically assess how effectively these models can use various tools, focusing on their proficiency with SQL and other coding languages.

The Effectiveness of simple SQL Creation Using the schemas provided in Table 12 and Table 13, we construct questions similar to those in Table 14, and refer readers to Appendix A. These questions are posed to various LLMs using our specifically designed prompts in Appendix B.

Following the tailored prompts, the LLMs are evaluated based on their responses to the presented queries. The results of this comprehensive assessment are compiled and exhibited in Figure 8.

This verifies the capabilities of each LLM in handling varying simple single-table SQL queries, thus providing a basis for comparison and analysis.

Table 8: The evaluation results for simple SQL questions.

Model	ChatGPT	Claude	Ziya
Accuracy	90%	100%	50%
Model	ChatGLM	Chinese-Alpaca-Plus	InternLM
Accuracy	30%	20%	90%

The Effectiveness of Complex Nested SQL Creation Using the schemas provided in Table 15, 16, 17, and 18, we construct questions similar to those in Table 19, and refer readers to Appendix A. For complex nested SQL questions, to further verify the SQL tool creation capability of LLMs, we have designed two types of prompts. One is the direct-guidance type, which explicitly informs the model that it needs to generate nested SQL query statements, as shown in Figure 14 in Appendix B.

The other is based on the Chain-of-Thought (CoT) [26] approach, which leverages the model’s ability to reason step by step to comprehend and craft SQL tools, and the prompt is shown in Figure 15 in Appendix B. This method guides the model to sequentially generate SQL query clauses based on the problem context, thus breaking down the complex query generation task into smaller and manageable subtasks. This approach provides the model with a structured way to handle complex SQL tasks and showcases its capacity to engage in incremental reasoning and problem-solving.

The design of these two types of prompts serves as the backbone of our evaluation for complex nested SQL questions. While the direct-guidance approach focuses on testing the model’s raw ability to generate SQL queries when explicitly instructed, the CoT-based approach evaluates a more nuanced capability: the model’s reasoning and problem-solving skills in a step-by-step manner. Both these methods present unique challenges and offer valuable insights into the strengths and potential areas of improvement for the large language model’s SQL tool generation ability. Subsequently, we will explore these two dimensions based on our experimental evaluations shown in Table 9.

Table 9: The evaluation results for complex nested SQL questions.

Model	ChatGPT	Claude	Ziya
Direct-based	80%	100%	50%
CoT-based	80%	100%	40%
Model	ChatGLM	Chinese-Alpaca-Plus	InternLM
Direct-based	60%	0%	60%
CoT-based	70%	0%	50%

From the above results in Table 9, it is clear that different models possess varying levels of proficiency in handling complex nested SQL tasks. Some models, like Claude, exhibit a robust capability in SQL generation, no matter whether the approach is direct or CoT-based. Most of these models demonstrate the SQL tool usage capability.

Specifically, some models such as ChatGLM show a distinct preference for the CoT-based approach, their performance improves when problems are broken down into smaller, manageable sub-tasks. This suggests that these models may have a stronger ability in sequential problem-solving and benefit more from step-by-step guidance. Conversely, models like Ziya and InternLM show a drop in performance when tasks are guided in the CoT-based format. This might indicate challenges in managing dependencies between sub-tasks or handling the continuity in sequential problem-solving. Lastly, Chinese-Alpaca-Plus shows significant room for improvement in complex SQL generation tasks. This shows that not all models are equally suited to handle advanced problem-solving involving nested SQL queries.

Overall, these findings underscore the importance of tailoring evaluation and training methodologies to the individual strengths and weaknesses of each model. By adopting this approach, we can better understand the performance variations across different models and provide targeted improvements to enhance their problem-solving abilities. Furthermore, this analysis highlights the potential of

LLM-based agents in real-world applications, and the need to push their boundaries through continued research and development.

The Effectiveness of Mathematical Code Creation Following our evaluation of the LLM’s proficiency in creating complex SQL queries, we now shift our focus to another tool creation: the creation of mathematical code. To the best of our knowledge, while large language models possess significant capabilities, they often fall short of providing highly accurate solutions to mathematical problems. Guiding these LLMs to generate mathematical code, and subsequently leveraging external tools to execute and derive the solutions, could significantly enhance their ability to tackle mathematical challenges.

In the upcoming section, we will conduct a detailed evaluation of guiding these LLMs to generate mathematical code. We aim to shed light on the true capability of these models in generating mathematical code and to elucidate the extent to which they can be utilized to aid in mathematical problem-solving. The prompt about how to guide LLMs is shown in Figure 16 in Appendix B.

Table 10: The evaluation results for mathematical questions.

Model	ChatGPT	Claude	Ziya
Accuracy	90%	85%	50%
Model	ChatGLM	Chinese-Alpaca-Plus	InternLM
Accuracy	0%	55%	95%

The results shown in Table 10 indicate that the capabilities of LLM-based agents to generate mathematical code vary considerably. High-performing models like ChatGPT, Claude, and InternLM display excellent proficiency, suggesting their potent ability to solve complex mathematical tasks. Middle-tier models, such as Ziya, show moderate success, indicating the potential for improvement and adaptability with the right training and optimization. Surprisingly, Alpaca demonstrated a notable proficiency in mathematical tasks, despite its poor performance in SQL generation, suggesting a possible inclination towards mathematical problems. In contrast, ChatGLM struggles significantly with mathematical code generation, underlining a potential weak spot in its capabilities and the need for focused improvement in this area.

Overall, these results underscore the task-dependent nature of LLMs’ capabilities and highlight the importance of recognizing their individual strengths and weaknesses for optimal model guidance and enhanced problem-solving.

3.3.2 TPTU-OA and TPTU-SA: Tool Usage for Multiple Tools

We now aim to utilize the one-step agent and sequential agent, which we designed, to conduct an evaluation involving multiple tools. Corresponding prompts for each agent type have been crafted and are presented in Figure 17 and Figure 18 of Appendix B, respectively.

In this phase of the evaluation, we need to automatically invoke the respective tools through code and produce the results. Given that user interface-based LLMs lack the capability to call external tools, we will only utilize the following four API-based LLMs (ChatGPT, Ziya, Chinese-Alpaca, and InternLM) for this comprehensive evaluation of external tool usage ability.

Table 11: The evaluation results for end-to-end ability of multiple tools.

Model	ChatGPT	Ziya	Chinese-Alpaca-Plus	InternLM
TPTU-OA	50%	0%	0%	15%
TPTU-SA	55%	0%	0%	20%

With agents mentioned above, the final results are presented in Table 11. The evaluation results demonstrate varying levels of task planning and tool usage capabilities among the four API-based LLMs. In the TPTU-OA evaluation, ChatGPT achieved a performance rate of 50%, significantly outperforming the other models, with InternLM at 15%, while both Ziya and Chinese-Alpaca did not manage to complete any tasks successfully, resulting in a score of 0%. In the TPTU-SA evaluation,

an overall slight improvement was observed. ChatGPT maintained its leading position, with a slightly improved performance rate of 55%. InternLM also exhibited better performance, achieving a score of 20%, whereas Ziya and Chinese-Alpaca-Plus again failed to register any successful task completion.

These results reflect a notable discrepancy in the performance of LLMs when it comes to using external tools. ChatGPT and InternLM have demonstrated some ability to navigate these tasks, but their performance rates suggest there is significant room for improvement. Ziya and Chinese-Alpaca-Plus' performance indicates a struggle to effectively utilize external tools in their current state.

The differential performance between the TPTU-OA and TPTU-SA evaluation hints at the possible impact of the agent design on the LLMs' task execution ability. In particular, the performance increase under the sequential agent framework suggests that breaking down tasks into sequential steps might help LLM-based AI agents better utilize external tools. This insight could prove valuable in future improvements and developments of LLM-based AI agents. However, even with this approach, it is clear that LLM-based AI agents are far from perfect when it comes to effectively using external tools for complex tasks. This finding underlines the importance of further investigation and improvement in this domain.

3.4 Insightful Observations

Upon closer observation of our experimental results, we have identified several phenomena that deserved further exploration. These findings serve to broaden our understanding of LLM-based agents' behavior and capabilities and provide essential insights that could shape future research in this field. In the following, we will dissect these phenomena as shown in Figure 4 - 7, casting light on the weaknesses of LLM-based agents in the context of task planning and tool usage.

- Misunderstanding Output Formats:** LLMs frequently encounter difficulty when output is required in specific formats such as lists or dictionaries. One such example includes inconsistencies between the number of tools and corresponding subtasks, leading to formatting issues that hinder the correct execution of tasks.

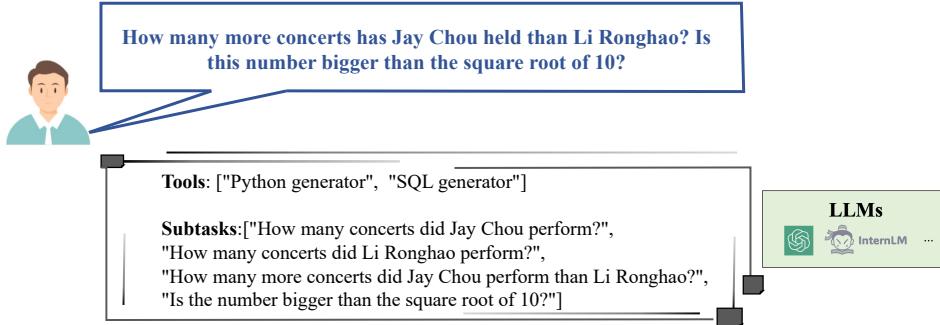


Figure 4: **Issue-1:** Inconsistencies between the number of tools and corresponding subtasks.

- Struggling to Grasp Task Requirements:** LLMs might incorrectly disintegrate subproblems or apply unsuitable tools to carry out the subproblem. For example, an LLM might attempt to solve a purely mathematical problem by employing an SQL tool or could misunderstand similar terms like cube extraction and cube roots.
- Endless Extensions:** LLMs tend to overutilize a particular tool, even in instances where a single use would suffice for the correct result. This issue can lead to extended and nonsensical planning, where the same subtask is repeatedly solved.
- Lack of Summary Skills:** LLMs do not take into account the responses to subproblems, relying instead on their internalized knowledge to generate the final answer. This may lead to a scenario where the final response only addresses a portion of the original query.

By identifying and addressing these common issues, we stand a better chance at improving and refining LLMs, thereby unlocking their full potential.

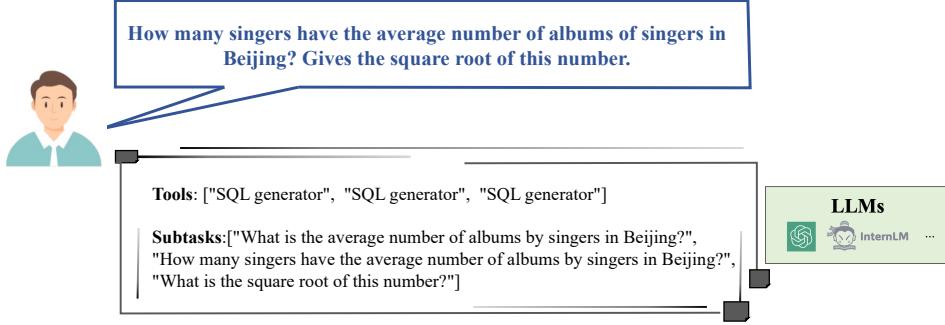


Figure 5: **Issue-2:**Solve a purely mathematical problem by employing a SQL generator.

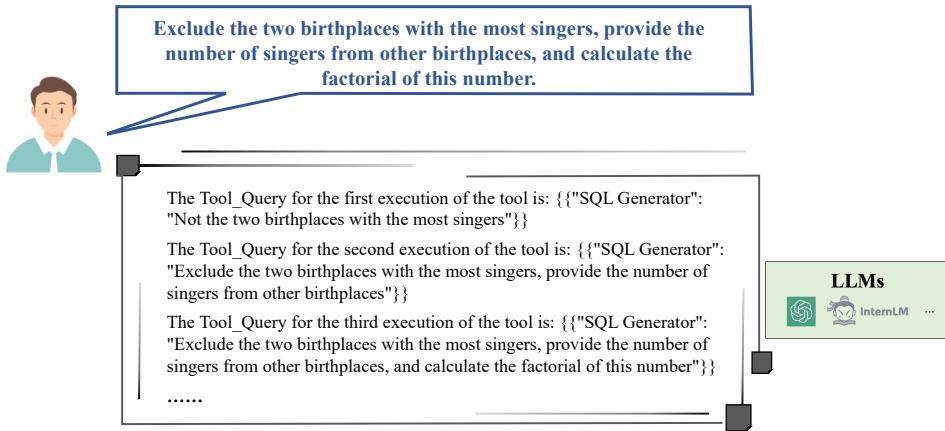


Figure 6: **Issue-3:** Unnecessary repetition of subtasks.

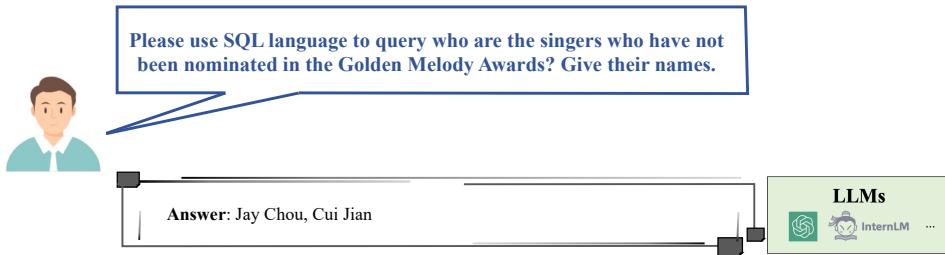


Figure 7: **Issue-4:** Answering questions using common sense instead of generating code.

4 Related Work

The remarkable capacity for usage and creation of tools have facilitated the transcendence of our innate physical and cognitive constraints, thereby profoundly advancing the progress and prosperity of human civilization and society. The swift advancement of LLM has rendered it feasible to use and create tools like humans. The integration of specialized tools with LLM has unlocked substantial potential in addressing intricate tasks. In this section, we offer a concise synopsis of the relevant research pertaining to tool learning based on LLMs.

4.1 Tool Usage

The initial advancements in tool learning have been constrained by the capabilities of artificial intelligence (AI) models. [27] Traditional deep learning approaches exhibit limitations in terms of comprehension of tool functionality and user intentions, and common sense reasoning abilities. Consequently, these limitations directly result in a notable decline in the stability and precision of tool

learning methodologies. Recently, the advent of LLM has marked a pivotal juncture in the realm of tool learning. LLMs encompass a broad spectrum of common sense cognitive capabilities and exhibit remarkable proficiencies in natural language processing, reasoning, and interactive decision-making [28–32]. These attributes furnish indispensable prerequisites for LLMs to comprehend user intentions and effectively employ tools in tackling intricate tasks [33]. Simultaneously, the advancement of fine-tuning [34–38] and in-context learning [39, 40] technology has offered robust support to LLM in addressing increasingly intricate challenges. In addition, tool usage can mitigate the inherent limitations of LLMs, encompassing the acquisition of up-to-date information from real-world events, refined mathematical computational abilities, and the mitigation of potential hallucinatory phenomena. [41]

Within the realm of embodied intelligence [42–44], LLM engages in direct interactions with tangible tools like robots in order to enhance their cognitive abilities, optimize work productivity, and expand functional capacities. LLM possesses the capability to automatically devise action steps based on user intentions, enabling the guidance of robots in the completion of tasks [45–53], or alternatively, to directly generate underlying code that can be executed by robots [54–58]. Palm-E [50] introduced a multimodal language model which seamlessly integrates sensor data into its framework, enabling efficient planning of robot actions and task completion. Code as Policies (CaP) [58] facilitates the transformation of natural language instructions into code fragments that can be directly compiled and executed on robots. As for Inner Monologue [48], LLM incorporates diverse environmental feedback to construct inner monologues, thereby formulating effective robot control strategies. Furthermore, LP-SLAM [45] proposes a simultaneous localization and mapping (SLAM) system empowered with language perception capabilities, exploiting the potential of ChatGPT. PromptCraft [57], on the other hand, devises a function library tailored to ChatGPT on the robot platform, streamlining the conversion of user intentions into executable tasks via the underlying backend API.

In addition to directly changing the real environment through interaction with tools in the physical world, LLM can also utilize software tools such as search engines [59–67], mobile [68, 69], Microsoft Office [70, 71], calculators [72–74], deep models [19, 75–79, 13, 80, 81] and other versatile APIs [82, 5, 83, 84, 20, 85] to enhance model performance or complete complex workflows through flexible control of the software. Toolformer [5] employs a self-supervised methodology to fine-tune the language model, enabling it to acquire the ability to automatically invoke APIs. ART [86] leverages CoT [26] and In-context Learning [81, 41] techniques to automatically generate multi-step reasoning processes for new tasks, while also selecting and utilizing the most appropriate available tool at each step. ASH [62] utilizes LLM for sequence hierarchical decision-making to achieve web navigation tasks. WebGPT [66] and WebCPM [64] use network search to assist in implementing Question Answering tasks. In addition, RCI [87] recursively criticizes and improves itself to execute computer tasks guided by natural language according to the prompting scheme. To achieve the analysis and processing of tables, TableGPT [71] employs a table encoder to transform tabular data into vector representations, which are then fed into an LLM for inference in combination with user queries.

4.2 Tool Creation

The usage of tools is contingent upon the accessibility of external tools. Recently, efforts have been made to employ LLM as a tool creator in order to generate tools that can be utilized for diverse requests [88–95]. This development has consequently raised the demands placed on LLM. And these created tools are typically implemented as Python or SQL functions. LATM [88], for example, leverages the prowess of GPT-4 to create tools, and the usage of more cost-effective models has shown potential in exhibiting performance on par with larger models for these tool applications. EVAPORATE [94] involves the synthesis of multiple functions, which are subsequently utilized at a large scale to efficiently process documents and generate structured views.

5 Conclusion

In this paper, we have introduced a structured framework specially designed for LLM-based AI Agents, with an emphasis on their abilities in task planning and tool usage. This framework, coupled with our design of two distinct types of agents assigned for the inference process, allows for a comprehensive evaluation of the capabilities of current open-source LLMs, thereby yielding critical insights into their effectiveness. Furthermore, our research highlights the significant potential of

LLMs in managing complex tasks, revealing the exciting prospects they hold for future research and development. As we continue to explore and improve upon these models, we move closer to unlocking their full potential in a wide range of real-world applications.

Acknowledgements

This work was conducted collaboratively among the authors.

Hangyu Mao and Rui Zhao led the project, formulating the central idea and laying out the framework for the primary literature review.

Regarding the literature review phase, the surveys were conducted by various team members. Guoqing Du and Jingqing Ruan explored DNN-based Tool Scheduling by LLMs; Tianpeng Bao and Yihong Chen investigated Physical/Robot Tool Scheduling by LLMs; and Shiwei Shi and Zhiwei Xu handled the survey of API or GUI-based Tool Scheduling by LLMs. Bin Zhang summarized these papers and synthesized an overarching summary.

As for the evaluation phase, Yihong Chen, Tianpeng Bao, Jingqing Ruan, Guoqing Du, Zhiwei Xu, Shiwei Shi, and Bin Zhang performed the experiments and analyzed the data. Hangyu Mao assisted in the analysis of the experimental phenomena and offered constructive suggestions for improvements. Xingyu Zeng and Rui Zhao provided invaluable feedback, contributed to the direction of the research. All authors participated in the discussion.

Regarding the manuscript phase, Hangyu Mao organized the overall chapters of the manuscript and mainly wrote the methodology part, and provided assistance in other parts. Jingqing Ruan and Yihong Chen wrote the evaluation section. Bin Zhang wrote the summary of the literature review. Each author read and approved the final manuscript.

The authors would like to thank Feng Zhu, Kun Wang, Yuhang Ran, Mengying Xu, Pengfei Jia, and Shaobo Lin for their valuable feedback, discussion, and participation in this project.

References

- [1] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, “A survey of large language models,” *arXiv preprint arXiv:2303.18223*, 2023.
- [2] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [3] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le, “Finetuned language models are zero-shot learners,” *arXiv preprint arXiv:2109.01652*, 2021.
- [4] OpenAI, “Gpt-4 technical report,” 2023.
- [5] T. Schick, J. Dwivedi-Yu, R. Dessì, R. Raileanu, M. Lomeli, L. Zettlemoyer, N. Cancedda, and T. Scialom, “Toolformer: Language models can teach themselves to use tools,” *arXiv preprint arXiv:2302.04761*, 2023.
- [6] N. R. Jennings, K. Sycara, and M. Wooldridge, “A roadmap of agent research and development,” *Autonomous agents and multi-agent systems*, vol. 1, pp. 7–38, 1998.
- [7] N. R. Jennings and M. Wooldridge, “Applying agent technology,” *Applied Artificial Intelligence an International Journal*, vol. 9, no. 4, pp. 357–369, 1995.
- [8] S. Franklin and A. Graesser, “Is it an agent, or just a program?: A taxonomy for autonomous agents,” in *International workshop on agent theories, architectures, and languages*. Springer, 1996, pp. 21–35.
- [9] C. Castelfranchi, “Modelling social action for ai agents,” *Artificial intelligence*, vol. 103, no. 1-2, pp. 157–182, 1998.
- [10] J. Ferber and G. Weiss, *Multi-agent systems: an introduction to distributed artificial intelligence*. Addison-wesley Reading, 1999, vol. 1.

- [11] L. Panait and S. Luke, “Cooperative multi-agent learning: The state of the art,” *Autonomous agents and multi-agent systems*, vol. 11, pp. 387–434, 2005.
- [12] M. Pourreza and D. Rafiei, “Din-sql: Decomposed in-context learning of text-to-sql with self-correction,” *arXiv preprint arXiv:2304.11015*, 2023.
- [13] C. Wu, S. Yin, W. Qi, X. Wang, Z. Tang, and N. Duan, “Visual chatgpt: Talking, drawing and editing with visual foundation models,” *arXiv preprint arXiv:2303.04671*, 2023.
- [14] J. Gorniak, Y. Kim, S. Gwon, D. Wei, and N. W. Kim, “Vizability: Multimodal accessible data visualization with keyboard navigation and conversational interaction,” *arXiv preprint arXiv:2310.09611*, 2023.
- [15] I. Team, “Internlm: A multilingual language model with progressively enhanced capabilities,” <https://github.com/InternLM/InternLM>, 2023.
- [16] Q. Xu, F. Hong, B. Li, C. Hu, Z. Chen, and J. Zhang, “On the tool manipulation capability of open-source large language models,” *arXiv preprint arXiv:2305.16504*, 2023.
- [17] Y. Qin, S. Liang, Y. Ye, K. Zhu, L. Yan, Y. Lu, Y. Lin, X. Cong, X. Tang, B. Qian *et al.*, “Toolllm: Facilitating large language models to master 16000+ real-world apis,” *arXiv preprint arXiv:2307.16789*, 2023.
- [18] M. Li, F. Song, B. Yu, H. Yu, Z. Li, F. Huang, and Y. Li, “Api-bank: A benchmark for tool-augmented llms,” *arXiv preprint arXiv:2304.08244*, 2023.
- [19] S. G. Patil, T. Zhang, X. Wang, and J. E. Gonzalez, “Gorilla: Large language model connected with massive apis,” *arXiv preprint arXiv:2305.15334*, 2023.
- [20] Q. Tang, Z. Deng, H. Lin, X. Han, Q. Liang, and L. Sun, “Toolalpaca: Generalized tool learning for language models with 3000 simulated cases,” *arXiv preprint arXiv:2306.05301*, 2023.
- [21] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, “Training language models to follow instructions with human feedback,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 27 730–27 744, 2022.
- [22] Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirhoseini, C. McKinnon *et al.*, “Constitutional ai: Harmlessness from ai feedback,” *arXiv preprint arXiv:2212.08073*, 2022.
- [23] A. Zeng, X. Liu, Z. Du, Z. Wang, H. Lai, M. Ding, Z. Yang, Y. Xu, W. Zheng, X. Xia *et al.*, “Glm-130b: An open bilingual pre-trained model,” *arXiv preprint arXiv:2210.02414*, 2022.
- [24] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, “Llama: Open and efficient foundation language models,” *arXiv preprint arXiv:2302.13971*, 2023.
- [25] Y. Cui, Z. Yang, and X. Yao, “Efficient and effective text encoding for chinese llama and alpaca,” *arXiv preprint arXiv:2304.08177*, 2023.
- [26] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le, and D. Zhou, “Chain-of-thought prompting elicits reasoning in large language models,” *Neural Information Processing Systems*, 2022.
- [27] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill *et al.*, “On the opportunities and risks of foundation models,” *arXiv preprint arXiv:2108.07258*, 2021.
- [28] M. Mosbach, T. Pimentel, S. Ravfogel, D. Klakow, and Y. Elazar, “Few-shot fine-tuning vs. in-context learning: A fair comparison and evaluation,” *arXiv preprint arXiv:2305.16938*, 2023.
- [29] J. Yang, H. Jin, R. Tang, X. Han, Q. Feng, H. Jiang, B. Yin, and X. Hu, “Harnessing the power of llms in practice: A survey on chatgpt and beyond,” *arXiv preprint arXiv:2304.13712*, 2023.

- [30] C. Zhang, C. Zhang, C. Li, Y. Qiao, S. Zheng, S. K. Dam, M. Zhang, J. U. Kim, S. T. Kim, J. Choi *et al.*, “One small step for generative ai, one giant leap for agi: A complete survey on chatgpt in aigc era,” *arXiv preprint arXiv:2304.06488*, 2023.
- [31] F. Yu, H. Zhang, and B. Wang, “Nature language reasoning, a survey,” *arXiv preprint arXiv:2303.14725*, 2023.
- [32] Z. Wang, G. Zhang, K. Yang, N. Shi, W. Zhou, S. Hao, G. Xiong, Y. Li, M. Y. Sim, X. Chen *et al.*, “Interactive natural language processing,” *arXiv preprint arXiv:2305.13246*, 2023.
- [33] Y. Qin, S. Hu, Y. Lin, W. Chen, N. Ding, G. Cui, Z. Zeng, Y. Huang, C. Xiao, C. Han *et al.*, “Tool learning with foundation models,” *arXiv preprint arXiv:2304.08354*, 2023.
- [34] W. Yu, C. Zhu, Z. Li, Z. Hu, Q. Wang, H. Ji, and M. Jiang, “A survey of knowledge-enhanced text generation,” *ACM Computing Surveys*, vol. 54, no. 11s, pp. 1–38, 2022.
- [35] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, “Lora: Low-rank adaptation of large language models,” *arXiv preprint arXiv:2106.09685*, 2021.
- [36] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly, “Parameter-efficient transfer learning for nlp,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 2790–2799.
- [37] X. L. Li and P. Liang, “Prefix-tuning: Optimizing continuous prompts for generation,” *arXiv preprint arXiv:2101.00190*, 2021.
- [38] X. Liu, Y. Zheng, Z. Du, M. Ding, Y. Qian, Z. Yang, and J. Tang, “Gpt understands, too,” *arXiv preprint arXiv:2103.10385*, 2021.
- [39] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao, “React: Synergizing reasoning and acting in language models,” *arXiv preprint arXiv:2210.03629*, 2022.
- [40] T. Khot, H. Trivedi, M. Finlayson, Y. Fu, K. Richardson, P. Clark, and A. Sabharwal, “Decomposed prompting: A modular approach for solving complex tasks,” *arXiv preprint arXiv:2210.02406*, 2022.
- [41] G. Mialon, R. Dessì, M. Lomeli, C. Nalmpantis, R. Pasunuru, R. Raileanu, B. Rozière, T. Schick, J. Dwivedi-Yu, A. Celikyilmaz *et al.*, “Augmented language models: a survey,” *arXiv preprint arXiv:2302.07842*, 2023.
- [42] J. Duan, S. Yu, H. L. Tan, H. Zhu, and C. Tan, “A survey of embodied ai: From simulators to research tasks,” *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 6, no. 2, pp. 230–244, 2022.
- [43] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik *et al.*, “Habitat: A platform for embodied ai research,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 9339–9347.
- [44] S. Franklin, “Autonomous agents as embodied ai,” *Cybernetics & Systems*, vol. 28, no. 6, pp. 499–520, 1997.
- [45] W. Zhang, Y. Guo, L. Niu, P. Li, C. Zhang, Z. Wan, J. Yan, F. U. D. Farrukh, and D. Zhang, “Lp-slam: Language-perceptive rgb-d slam system based on large language model,” *arXiv preprint arXiv:2303.10089*, 2023.
- [46] D. Shah, B. Osiński, S. Levine *et al.*, “Lm-nav: Robotic navigation with large pre-trained models of language, vision, and action,” in *Conference on Robot Learning*. PMLR, 2023, pp. 492–504.
- [47] A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang, R. Julian *et al.*, “Do as i can, not as i say: Grounding language in robotic affordances,” in *Conference on Robot Learning*. PMLR, 2023, pp. 287–318.

- [48] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mordatch, Y. Chebotar *et al.*, “Inner monologue: Embodied reasoning through planning with language models,” *arXiv preprint arXiv:2207.05608*, 2022.
- [49] B. Chen, F. Xia, B. Ichter, K. Rao, K. Gopalakrishnan, M. S. Ryoo, A. Stone, and D. Kappler, “Open-vocabulary queryable scene representations for real world planning,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 11509–11522.
- [50] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu *et al.*, “Palm-e: An embodied multimodal language model,” *arXiv preprint arXiv:2303.03378*, 2023.
- [51] N. Wake, A. Kanehira, K. Sasabuchi, J. Takamatsu, and K. Ikeuchi, “Chatgpt empowered long-step robot control in various environments: A case application,” *arXiv preprint arXiv:2304.03893*, 2023.
- [52] K. Rana, J. Haviland, S. Garg, J. Abou-Chakra, I. Reid, and N. Suenderhauf, “Sayplan: Grounding large language models using 3d scene graphs for scalable task planning,” *arXiv preprint arXiv:2307.06135*, 2023.
- [53] C. H. Song, J. Wu, C. Washington, B. M. Sadler, W.-L. Chao, and Y. Su, “Llm-planner: Few-shot grounded planning for embodied agents with large language models,” *arXiv preprint arXiv:2212.04088*, 2022.
- [54] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu *et al.*, “Rt-1: Robotics transformer for real-world control at scale,” *arXiv preprint arXiv:2212.06817*, 2022.
- [55] A. Stone, T. Xiao, Y. Lu, K. Gopalakrishnan, K.-H. Lee, Q. Vuong, P. Wohlhart, B. Zitkovich, F. Xia, C. Finn *et al.*, “Open-world object manipulation using pre-trained vision-language models,” *arXiv preprint arXiv:2303.00905*, 2023.
- [56] S. Reed, K. Zolna, E. Parisotto, S. G. Colmenarejo, A. Novikov, G. Barth-Maron, M. Gimenez, Y. Sulsky, J. Kay, J. T. Springenberg *et al.*, “A generalist agent,” *arXiv preprint arXiv:2205.06175*, 2022.
- [57] S. Vemprala, R. Bonatti, A. Bucker, and A. Kapoor, “Chatgpt for robotics: Design principles and model abilities,” *Microsoft Auton. Syst. Robot. Res.*, vol. 2, p. 20, 2023.
- [58] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng, “Code as policies: Language model programs for embodied control,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 9493–9500.
- [59] K. Guu, K. Lee, Z. Tung, P. Pasupat, and M. Chang, “Retrieval augmented language model pre-training,” in *International conference on machine learning*. PMLR, 2020, pp. 3929–3938.
- [60] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel *et al.*, “Retrieval-augmented generation for knowledge-intensive nlp tasks,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 9459–9474, 2020.
- [61] S. Borgeaud, A. Mensch, J. Hoffmann, T. Cai, E. Rutherford, K. Millican, G. B. Van Den Driessche, J.-B. Lespiau, B. Damoc, A. Clark *et al.*, “Improving language models by retrieving from trillions of tokens,” in *International conference on machine learning*. PMLR, 2022, pp. 2206–2240.
- [62] A. Sridhar, R. Lo, F. F. Xu, H. Zhu, and S. Zhou, “Hierarchical prompting assists large language model on web navigation,” *arXiv preprint arXiv:2305.14257*, 2023.
- [63] H. Furuta, O. Nachum, K.-H. Lee, Y. Matsuo, S. S. Gu, and I. Gur, “Multimodal web navigation with instruction-finetuned foundation models,” *arXiv preprint arXiv:2305.11854*, 2023.
- [64] Y. Qin, Z. Cai, D. Jin, L. Yan, S. Liang, K. Zhu, Y. Lin, X. Han, N. Ding, H. Wang *et al.*, “Webcpm: Interactive web search for chinese long-form question answering,” *arXiv preprint arXiv:2305.06849*, 2023.

- [65] S. Yao, H. Chen, J. Yang, and K. Narasimhan, “Webshop: Towards scalable real-world web interaction with grounded language agents,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 20744–20757, 2022.
- [66] R. Nakano, J. Hilton, S. Balaji, J. Wu, L. Ouyang, C. Kim, C. Hesse, S. Jain, V. Kosaraju, W. Saunders *et al.*, “Webgpt: Browser-assisted question-answering with human feedback,” *arXiv preprint arXiv:2112.09332*, 2021.
- [67] Z. Yang, P. Qi, S. Zhang, Y. Bengio, W. W. Cohen, R. Salakhutdinov, and C. D. Manning, “Hotpotqa: A dataset for diverse, explainable multi-hop question answering,” *arXiv preprint arXiv:1809.09600*, 2018.
- [68] B. Wang, G. Li, and Y. Li, “Enabling conversational interaction with mobile ui using large language models,” in *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, 2023, pp. 1–17.
- [69] D. Zhang, L. Chen, and K. Yu, “Mobile-env: A universal platform for training and evaluation of mobile interaction,” *arXiv preprint arXiv:2305.08144*, 2023.
- [70] H. Li, J. Su, Y. Chen, Q. Li, and Z. Zhang, “Sheetcopilot: Bringing software productivity to the next level through large language models,” *arXiv preprint arXiv:2305.19308*, 2023.
- [71] L. Zha, J. Zhou, L. Li, R. Wang, Q. Huang, S. Yang, J. Yuan, C. Su, X. Li, A. Su *et al.*, “Tablegpt: Towards unifying tables, nature language and commands into one gpt,” *arXiv preprint arXiv:2307.08674*, 2023.
- [72] Z. Chen, K. Zhou, B. Zhang, Z. Gong, W. X. Zhao, and J.-R. Wen, “Chatcot: Tool-augmented chain-of-thought reasoning on \chat-based large language models,” *arXiv preprint arXiv:2305.14323*, 2023.
- [73] A. Parisi, Y. Zhao, and N. Fiedel, “Talm: Tool augmented language models,” *arXiv preprint arXiv:2205.12255*, 2022.
- [74] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano *et al.*, “Training verifiers to solve math word problems,” *arXiv preprint arXiv:2110.14168*, 2021.
- [75] Z. Yang, L. Li, J. Wang, K. Lin, E. Azarnasab, F. Ahmed, Z. Liu, C. Liu, M. Zeng, and L. Wang, “Mm-react: Prompting chatgpt for multimodal reasoning and action,” *arXiv preprint arXiv:2303.11381*, 2023.
- [76] Z. Liu, Y. He, W. Wang, Y. Wang, S. Chen, Q. Zhang, Y. Yang, Q. Li, J. Yu *et al.*, “Internchat: Solving vision-centric tasks by interacting with chatbots beyond language,” *arXiv preprint arXiv:2305.05662*, 2023.
- [77] Y. Ge, W. Hua, J. Ji, J. Tan, S. Xu, and Y. Zhang, “Openagi: When llm meets domain experts,” *arXiv preprint arXiv:2304.04370*, 2023.
- [78] Y. Shen, K. Song, X. Tan, D. Li, W. Lu, and Y. Zhuang, “Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface,” *arXiv preprint arXiv:2303.17580*, 2023.
- [79] D. Surís, S. Menon, and C. Vondrick, “Vipergpt: Visual inference via python execution for reasoning,” *arXiv preprint arXiv:2303.08128*, 2023.
- [80] T. Gupta and A. Kembhavi, “Visual programming: Compositional visual reasoning without training,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14953–14962.
- [81] L. Chen, B. Li, S. Shen, J. Yang, C. Li, K. Keutzer, T. Darrell, and Z. Liu, “Language models are visual reasoning coordinators,” in *ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2023.
- [82] P. Lu, B. Peng, H. Cheng, M. Galley, K.-W. Chang, Y. N. Wu, S.-C. Zhu, and J. Gao, “Chameleon: Plug-and-play compositional reasoning with large language models,” *arXiv preprint arXiv:2304.09842*, 2023.

- [83] Z. Gou, Z. Shao, Y. Gong, Y. Shen, Y. Yang, N. Duan, and W. Chen, “Critic: Large language models can self-correct with tool-interactive critiquing,” *arXiv preprint arXiv:2305.11738*, 2023.
- [84] Y. Liang, C. Wu, T. Song, W. Wu, Y. Xia, Y. Liu, Y. Ou, S. Lu, L. Ji, S. Mao *et al.*, “Taskmatrix.ai: Completing tasks by connecting foundation models with millions of apis,” *arXiv preprint arXiv:2303.16434*, 2023.
- [85] S. Hao, T. Liu, Z. Wang, and Z. Hu, “Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings,” *arXiv preprint arXiv:2305.11554*, 2023.
- [86] B. Paranjape, S. Lundberg, S. Singh, H. Hajishirzi, L. Zettlemoyer, and M. T. Ribeiro, “Art: Automatic multi-step reasoning and tool-use for large language models,” *arXiv preprint arXiv:2303.09014*, 2023.
- [87] G. Kim, P. Baldi, and S. McAleer, “Language models can solve computer tasks,” *arXiv preprint arXiv:2303.17491*, 2023.
- [88] T. Cai, X. Wang, T. Ma, X. Chen, and D. Zhou, “Large language models as tool makers,” *arXiv preprint arXiv:2305.17126*, 2023.
- [89] R. H. Lewis and J. Jiao, “Computeegpt: A computational chat model for numerical problems,” *arXiv preprint arXiv:2305.06223*, 2023.
- [90] L. Gao, A. Madaan, S. Zhou, U. Alon, P. Liu, Y. Yang, J. Callan, and G. Neubig, “Pal: Program-aided language models,” in *International Conference on Machine Learning*. PMLR, 2023, pp. 10 764–10 799.
- [91] G. Wang, Y. Xie, Y. Jiang, A. Mandlekar, C. Xiao, Y. Zhu, L. Fan, and A. Anandkumar, “Voyager: An open-ended embodied agent with large language models,” *arXiv preprint arXiv:2305.16291*, 2023.
- [92] C. Qian, C. Han, Y. R. Fung, Y. Qin, Z. Liu, and H. Ji, “Creator: Disentangling abstract and concrete reasonings of large language models through tool creation,” *arXiv preprint arXiv:2305.14318*, 2023.
- [93] Y. Cai, S. Mao, W. Wu, Z. Wang, Y. Liang, T. Ge, C. Wu, W. You, T. Song, Y. Xia *et al.*, “Low-code llm: Visual programming over llms,” *arXiv preprint arXiv:2304.08103*, 2023.
- [94] S. Arora, B. Yang, S. Eyuboglu, A. Narayan, A. Hojel, I. Trummer, and C. Ré, “Language models enable simple systems for generating structured views of heterogeneous data lakes,” *arXiv preprint arXiv:2304.09433*, 2023.
- [95] W. Zhang, Y. Shen, W. Lu, and Y. Zhuang, “Data-copilot: Bridging billions of data and humans with autonomous workflow,” *arXiv preprint arXiv:2306.07209*, 2023.

A Detailed Dataset Description

Simple SQL queries: These queries typically involve basic operations such as SELECT, FROM, WHERE, GROUP BY, etc. They are used to retrieve, filter, group, and sort data from a single table. We give the Schema of two tables in the SQL database in Table 12 and 13 and list several examples in Table 14.

Table 12: Schema of the Person table

Person	
Column Name	Type
id	TEXT
name	TEXT
age	INTEGER
sex	TEXT
school	TEXT
phone	TEXT
qualifications	TEXT
ability	TEXT

Table 13: Schema of the School table

School	
Column Name	Type
id	TEXT
name	TEXT
info_985	TEXT
info_211	TEXT

Table 14: Demonstrations of simple SQL queries.

Table ID	Question	Answer	SQL reference
Person	Average ages	35.16	select avg(age) from Person
Person	How many men	12	select count(*) from Person where sex = 'male'
School	How many schools are both '985' and '211' institutions?	11	select count(*) from School where info_985 = 'yes' and info_211 = 'yes';

Complex nested SQL queries: These queries contain subqueries, which are SQL queries nested inside a larger query. Nested queries can be used in various clauses such as SELECT, FROM, WHERE, and HAVING. They provide a way to perform multiple operations or calculations across multiple tables. We give the Schema of two tables in the SQL database in Table 15, 16, 17, and 18 and list several examples in Table 19.

Table 15: Schema of GoldenMelodyAwards

GoldenMelodyAwards	
Column Name	Type
Nominated_Count	INTEGER
Competing_Count	INTEGER
Awards_Count	INTEGER
Award_Name	TEXT
Host	TEXT
Year	TIME

Table 16: Schema of the AwardNominees table

AwardNominees	
Column Name	Type
Singer_ID	INTEGER
Nominated_Work	TEXT
Award_Name	TEXT
Award_Edition_ID	INTEGER

Complex nested queries utilizing multiple tools: These are advanced queries that involve multiple tools, such as SQL queries, python code generation, user-defined functions, etc. We give the Schema

Table 17: Schema of the Singers table

Singers	
Column Name	Type
Name	TEXT
Song_Count	INTEGER
Album_Count	INTEGER
Fan_Count	INTEGER
Gender	TEXT
Singer_ID	INTEGER

Table 18: Schema of the RecordCompanies table

RecordCompanies	
Column Name	Type
Record_Company	TEXT
Signing_Date	TIME
Singer_ID	INTEGER

Table 19: Demonstrations of complex nested SQL queries.

Table ID	Question	Answer	SQL reference
GoldenMelody-Awards	Golden Melody hosts, excluding the two with the least awards.	"26th Golden Melody", "27th Golden Melody"	select Award_Name from GoldenMelodyAwards where Host not in (select Host from GoldenMelodyAwards group by Host order by avg (Awards_Count) asc limit 2)
AwardNominees & Singers	Names of singers never nominated for Golden Melody Awards.	"Jay Chou", "Jian Cui"	select Name from Singers where Singer_ID not in (select Singer_ID from AwardNominees)
RecordCompanies & Singers	Name and gender of singers without a record company.	"Penny Tai:Femal"	select Name, Gender from Singers where Singer_ID not in (select Singer_ID from RecordCompanies);
GoldenMelody-Awards	How many times is the 27th Golden Melody count of the 28th's?	1	select a.Awards_Count / b.Awards_Count from (select Awards_Count from GoldenMelodyAwards where Award_Name == '27th Golden Melody') a , (select Awards_Count from GoldenMelodyAwards where Award_Name == '28th Golden Melody') b

of two tables in the SQL database in Table 20, and 21 and list several examples in Table 22. For verifying the planning ability of the LLM-based AI agents, we select this type of query.

Table 20: Schema of the Journal table

Journal	
Column Name	Type
Name	TEXT
First_Issue_Date	TIME
Journal_ID	INTEGER
Category	TEXT
Sponsor_Organization	TEXT
Country	TEXT
s Language	TEXT
Publication_Count	INTEGER

Table 21: Schema of the CoverPersonality table

CoverPersonality	
Column Name	Type
Person_ID	INTEGER
Journal_ID	INTEGER
Count	INTEGER

Table 22: Demonstrations of complex nested queries utilizing multiple tools.

Table ID	Question	Answer	Planning Tools	SQL reference	Code reference
Journal & Cover-Personality	Calculate the exponential of 3 and list the names and languages of journals with no cover personality.	[20.08, "The Economist: Chinese, Reader's Digest: English."]	["PythonREPL", "SQL Generator"]	select Name, Language from Journal where Journal_ID not in (select Journal_ID from CoverPersonality)	import math; return math.exp(3)
CoverPersonality & Journal	Compute 4's factorial, compare with GCD of 212 and list the names and languages of journals with no cover personality.	[4, "The Economist: Chinese, Reader's Digest: English."]	["PythonREPL", "SQL Generator"]	select Name, Language from Journal where Journal_ID not in (select Journal_ID from CoverPersonality)	import math; math.factorial(4), 212
Journal	Calculate the square root of 24, and query for the language whose average number of published issues exceeds the overall average.	[4.8989795, "English"]	["PythonREPL", "SQL Generator"]	select Language from Journal group by Language having avg (Publication_Count) > (select avg (Publication_Count) from Journal)	import math; math.sqrt(24)
CoverPersonality	Compute the log base 10 of 5, then identify cover figures appearing less than the overall max frequency across journals.	[0.69897, "Qing Hai, Xiaoming Huang, Cristiano Ronaldo, Kobe Bryant"]	["PythonREPL", "SQL Generator"]	select Person_ID from CoverPersonality where Count < (select max (Count) from CoverPersonality)	import math; math.log10(5)

B Prompts Design

Figure 8: The evaluation prompt for tool order planning.

You are a strategy model. Given a problem **and** a **set** of tools, you need to
↪ generate a sequence of tools to determine the solution to the problem.

Each tool **in** the toolset **is** defined **as** follows:

SQL Generator: Given an **input** problem **and** a database, it creates a
↪ syntactically correct SQLite query statement.

Python Generator: Given an **input** problem **and** some information, it generates
↪ a syntactically correct Python code snippet.

Please use the following **format**:

Question: This **is** the original question.

Error: This **is** the previously generated error output.

Tool: These are the tools to be selected **and** the order **in** which they are
↪ called. Please note to generate a Tool different **from** the Error.

Result: The final result output by the tool.

Here are some examples of mapping problems to tools:

Question: What **is** the square of the number of albums by Jolin Tsai?

Error: **None**

Tool: ["SQL Generator", "Python Generator"]

Result: 100

Question: First, calculate the square of 40, denoted **as** A, **and** then find
↪ the names of **all** the singers whose total number of fans **is** less than A.

Error: **None**

Tool: ["Python Generator", "SQL Generator"]

Result: ['Jolin Tsai']

Let's get started:

Question: {question}

Error: {error}

Tool:

Figure 9: The evaluation prompt for tool order and subtask description planning.

You are a strategy model. Given a problem **and** a **set** of tools, you need to
→ generate a sequence of tools to determine the solution to the problem.

Each tool **in** the toolset **is** defined **as** follows:

SQL Generator: Given an **input** problem **and** a database, it creates a
→ syntactically correct SQLite query statement.

Python Generator: Given an **input** problem **and** some information, it generates
→ a syntactically correct Python code snippet.

Please use the following **format**:

Question: This **is** the original question.

Error: This **is** the previously generated error output.

Tool: These are the tools to be selected **and** the order **in** which they are
→ called. Please note to generate a Tool different **from** the Error.

Query: This **is** the sub-problem derived **from** the original question that
→ needs to be **input** when calling the tool. Please note to generate a
→ Query different **from** the Error.

Result: The final result output by the tool.

Here are some examples of mapping problems to tools:

Question: What **is** the square of the number of albums by Jolin Tsai?

Error: **None**

Tool: ["SQL Generator", "Python Generator"]

Query: ["What is the number of albums by Jolin Tsai?", "What is the square
→ of the number of albums by Jolin Tsai?"]

Result: 100

Question: First, calculate the square of 40, denoted **as** A, **and** then find
→ the names of **all** the singers whose total number of fans **is** less than A.

Error: **None**

Tool: ["Python Generator", "SQL Generator"]

Query: ["A is the square of 40, what is the value of A?", "What are the
→ names of all the singers whose total number of fans is less than A?"]

Result: ['Jolin Tsai']

Let's get started:

Question: {question}

Error: {error}

Tool: {tools}

Query:

Figure 10: The evaluation prompt for one-step tool-subtask pair planning.

You are a strategy model. Given a problem **and** a **set** of tools, you need to
→ generate a sequence of tools to determine the solution to the problem.

Each tool **in** the toolset **is** defined **as** follows:

SQL Generator: Given an **input** problem **and** a database, it creates a
→ syntactically correct SQLite query statement.

Python Generator: Given an **input** problem **and** some information, it generates
→ a syntactically correct Python code snippet.

Please use the following **format**:

Question: This **is** the original question

Error: This **is** the previously generated error output

Tasks: This **is** a **list** **in** Python. Each item **in** the **list** **is** a dictionary. The
→ key of the dictionary represents the selected Tool, **and** the value **is**
→ the Query when calling the tool. Please note to generate a Tool **and**
→ Query different **from** the Error.

Answer: The final answer

Here are some examples of mapping problems to tools:

Question: What **is** the square of the number of albums by Jolin Tsai?

Error: **None**

Tasks: [{"SQL Generator": "What is the number of albums by Jolin Tsai?"},
→ {"Python Generator": "What is the square of the number of albums by
→ Jolin Tsai?"}]

Answer: The square of the number of albums by Jolin Tsai **is** 100

Question: First, calculate the square of 40, denoted **as** A, **and** then find
→ the names of **all** the singers whose total number of fans **is** less than A.

Error: **None**

Tasks: [{"Python Generator": "A is the square of 40, what is the value of
→ A?"}, {"SQL Generator": "What are the names of all the singers whose
→ total number of fans is less than A?"}]

Answer: Jolin Tsai

You must note that: The generated Tasks must strictly meet the **format**
→ requirements: it must be a **list** **in** Python, each item **in** the **list** **is** a
→ dictionary, the key of the dictionary represents the selected Tool, **and**
→ the value **is** the Query when calling the tool.

Let's get started:

Question: {question}

Error: {error}

Tasks: """

Figure 11: The prompt added to Figure 10 for tool-subtask pair planning with other unrelated tools.

Each tool `in` the toolset `is` defined `as` follows:

- SQL Generator: Given an `input` problem `and` a database, it creates a
 - ↪ syntactically correct SQLite query statement.
- Python Generator: Given an `input` problem `and` some information, it generates
 - ↪ a syntactically correct Python code snippet.
- Weather Query Tool: Given a location, it outputs the real-time weather of
 - ↪ that location.
- Image Generator: Given a text description, it generates a related image.
- Text Extractor: Given a link to an image, it extracts the corresponding
 - ↪ text `and` its position coordinates.
- Translator: Given a piece of text, it translates it into other languages.
- Bing Searcher: Given a piece of text, it conducts a search `in` the Bing
 - ↪ browser `and` returns the content.
- Shell Generator: Given an `input` problem `and` some information, it generates
 - ↪ a syntactically correct Shell script.
- Java Generator: Given an `input` problem `and` some information, it generates a
 - ↪ syntactically correct Java code snippet.
- Wikipedia Searcher: Given a piece of text, it conducts a search `in`
 - ↪ Wikipedia `and` returns the content.
- Office Suite: Given a text description, it automatically generates the
 - ↪ corresponding long document, table, `or` PPT.
- Movie Player: Given a movie name, it automatically plays the corresponding
 - ↪ movie resource.

Figure 12: The prompt for the tool-subtask pair generation with TPTU-SA.

```
You are a strategic model. Given a problem and a set of tools, you need to generate
↪ the next tool to be called and the corresponding subtask.

Each tool in the toolset is defined as follows:
SQL Generator: Given an input question and a database, it creates a syntactically
↪ correct SQLite query statement.
PythonREPL: Given an input question and some information, it generates a segment of
↪ syntactically correct Python code.

Please use the following format:

Question: This is the question
History: This is the history of previously generated sub-problems; if it's empty, it
↪ means there are no historical information currently
Tool_Query: This is a dictionary in Python, where the key represents the chosen
↪ Tool, and the value is the query input when invoking the Tool.
Result: This is the output result of the current Tool_Query Tool
...
History: This is the history of all previously generated sub-problems
Tool_Query: 'None' signifies that the Final_Answer can be derived
Result: 'None' signifies that the Final_Answer can be derived
Final_Answer: This is the final answer; when the history is sufficient to reason out
↪ the answer, provide the Final_Answer directly

In the above format, ... signifies that (History/Tool_Query/Result) can be repeated
↪ N times.
When you can get the Final_Answer, you can generate an empty Tool_Query and Result,
↪ and provide the Final_Answer
Please stop after generating the Result line or the Final_Answer line.

Below are some examples:

Question: First calculate the square of 40 as A, and find the names of all singers
↪ whose total fan count is less than A.
History:
Tool_Query:{{"PythonREPL": "A is the square of 40, what is the value of A?"}}
Result:1600
History: The Tool_Query for the first tool execution was:{{"PythonREPL": "A is the
↪ square of 40, what is the value of A?"}}, Result:1600
Tool_Query:{{"SQL Generator": "Find the names of all singers whose total fan count
↪ is less than A"}}
Result: Jolin Tsai

History: The Tool_Query for the first tool execution was: {{"PythonREPL": "A is the
↪ square of 40, what is the value of A?"}}, Result: 1600
      The Tool_Query for the second tool execution was: {{"SQL Generator": "Find
↪ the names of all singers whose total fan count is less than A"}},  

↪ Result: Jolin Tsai
Tool_Query:None
Result:None
Final_Answer: Jolin Tsai

Note: The generated Tool_Query must strictly comply with the format requirements,
↪ and only one Tool_Query can be generated each time. Do not perform additional
↪ problem analysis, strictly adhere to the format of the problem, and generate
↪ output similar to the examples.

Now let's get started:
Question: {question}
History: {history}
Tool_Query:
```

Figure 13: The evaluation prompt for simple SQL questions.

You are an SQLite expert. Given an `input` question, first, generate a
→ grammatically correct SQLite query to execute. Then examine the query
→ results and provide an answer to the `input` question.
Unless a specific number of examples to retrieve `is` specified `in` the
→ question, use the `LIMIT` clause to query `for` a maximum of 5 results.
Do `not` query `all` columns `in` the table. You must only query the columns
→ necessary to answer the question.
Please only use the column names you can see `in` the table below. Be careful
→ `not` to query columns that do `not` exist. Additionally, be aware of which
→ column `is in` which table.

Please use the following `format`:

Question: This `is` the question.
SQLQuery: The SQL query to be executed.
SQLResult: The result of the SQL query execution.
Answer: The final answer.

Note to only use the tables below:

```
CREATE TABLE Person (\n\t id TEXT, \n\t name TEXT, \n\t age INTEGER, \n\t\n\t sex TEXT, \n\t school TEXT, \n\t phone TEXT, \n\t qualifications TEXT,\n\t \n\t ability TEXT\n)\n/*\n 3 rows from person table:\n  id\tname\tage\tsex\tschool\tphone\tqualifications\tability\n  01\tWang Min\t32\tFemale\tBeijing University of Technology\t13938493271\tUndergraduate\tTourism Industry-related Work\n  02\tLi Liang\t27\tMale\tBeijing University of Technology\t13812764851\tMaster\tInternet Company Operations\n  03\tZhang Jing\t50\tFemale\tWuhan University of Technology\t13764592384\tMaster\tEditor of Publishing House\n*/\n\nCREATE TABLE School (\n\t id TEXT, \n\t name TEXT, \n\t info\_985 TEXT,\n\t \n\t info\_211 TEXT\n)\n/*\n 3 rows from school table:\n  id\tname\tinfo\_985\tinfo\_211\n  01\tCentral South University\tyes\tyes\n  02\tShandong University\tyes\tyes\n  03\tTsinghua University\tyes\tyes\n*/
```

Question: What `is` the average age of the people?
SQLQuery:

Figure 14: The evaluation prompt for complex nested SQL questions.

You are an SQL expert. Given an `input` question, you need to create a
→ syntactically correct SQL query statement. Please only use the
→ following datasets, which include four table names: `GoldenMelodyAward`,
→ `Singers`, `AwardNominee`, `Singers`, and `RecordCompanies`. The column names
→ and types of each table can be obtained from the create commands in the
→ table below:

```
CREATE TABLE GoldenMelodyAward (\n\t Nominated\_Count INTEGER, \n\t  
→ Competing\_Count INTEGER, \n\t Awards\_Count INTEGER, \n\t Award\_Name  
→ TEXT, \n\t Host TEXT, \n\t Year TIME \n) \n\n\nCREATE TABLE AwardNominees (\n\t Singer_ID INTEGER, \n\t Nominated\_Work  
→ TEXT, \n\t Award\_Name TEXT, \n\t Award_Edition_ID INTEGER \n) \n\n\nCREATE TABLE Singers(\n\t Name TEXT, \n\t Song\_Count INTEGER, \n\t  
→ Album\_Count INTEGER, \n\t Fan\_Count INTEGER, \n\t Singer\_ID INTEGER,  
→ \n\t Gender TEXT \n) \n\n\nCREATE TABLE RecordCompanies (\n\t Record\_Company TEXT, \n\t Singer\_Date  
→ TIME, \n\t Singer_ID INTEGER \n) \n\n\n
```

You can query one or more tables at the same time. Be careful not to query
→ non-existent table names or column names. Also, please note which
→ column is in which table.

Please use the following `format` when answering:

Question: This is the question

Answer: The SQL query statement to be executed

Figure 15: The evaluation CoT-based prompt for complex nested SQL questions.

You are an SQL expert. Given an `input` question, you need to create a
→ syntactically correct SQL query statement. Please only use the
→ following datasets, which include four table names: `GoldenMelodyAward`,
→ `Singers`, `AwardNominee`, `Singers`, and `RecordCompanie`. The column names
→ and types of each table can be obtained from the create commands in the
→ table below:

```
CREATE TABLE GoldenMelodyAward (\n\t Nominated\_Count INTEGER, \n\t  
→ Competing\_Count INTEGER, \n\t Awards\_Count INTEGER, \n\t Award\_Name  
→ TEXT, \n\t Host TEXT, \n\t Year TIME \n) \n\nCREATE TABLE AwardNominees (\n\t Singer_ID INTEGER, \n\t Nominated\_Work  
→ TEXT, \n\t Award\_Name TEXT, \n\t Award_Edition_ID INTEGER \n) \n\nCREATE TABLE Singers(\n\t Name TEXT, \n\t Song\_Count INTEGER, \n\t  
→ Album\_Count INTEGER, \n\t Fan\_Count INTEGER, \n\t Singer\_ID INTEGER,  
→ \n\t Gender TEXT \n) \n\nCREATE TABLE RecordCompanies (\n\t Record\_Company TEXT, \n\t Singer\_Date  
→ TIME, \n\t Singer_ID INTEGER \n) \n\n
```

You can query one or more tables at the same time. Be careful not to query
→ non-existent table names or column names. Also, please note which
→ column is in which table.

Please note that you are not proficient in nested SQL, when encountering
→ complex problems, you can think step by step to generate multiple
→ non-nested SQL statements. For example:

Question: Some minor languages are used by no more than 3 countries, what
→ are the source countries of these languages?
Thought: First generate the 1st SQL `select Official_Language from Country`
→ group by Official_Language having count(*) > 3', and assume that the
→ result of this SQL is result1, then generate the 2nd SQL `select Name`
→ `from Country where Official_Language not in result1'`.
Answer: `select Name from Country where Official_Language not in (select`
→ `Official_Language from Country group by Official_Language having`
→ `count(*) > 3)`

Please use the following format when answering:

Question: This is the question
Thought: This is the thought process
Answer: This is the final SQL query statement

Figure 16: The evaluation prompt for mathematical questions.

Transform a math problem into a solution function that can be executed using
→ Python's math library. Use the output of running this code to answer the
→ question.

Please use the following **format**:

History: Information output from previous tool invocation
Question: A question about mathematics
Error: This is the error output previously generated
PythonSolution: A Python solution, make sure to generate a PythonSolution different
→ from the one in Error, for example,
Python Solution
def solution():
 Python statement
Answer: The final answer

Below are some demonstrations of mapping math problems to PythonSolution:

History: The original question was: What is $37593 * 67$?
Question: What is $37593 * 67$?
Error: None
PythonSolution:
Python Solution
def solution():
 import math
 return 37593 * 67
Answer: 2518731

History: The original question was: What is the 1/5th power of 37593 ?
Question: What is the 1/5th power of 37593 ?
Error: None
PythonSolution:
Python Solution
def solution():
 import math
 return 37593 ** 1/5
Answer: 8.222831614237718

History: The original question was: What is the logarithm of 5 with base 10?
Question: What is the logarithm of 5 with base 10?
Error: None
PythonSolution:
Python Solution
def solution():
 import math
 return math.log(5, 10)
Answer: 0.69897

Now let's get started:

History:{history}
Question: {question}
Error:{error}
PythonSolution:

Figure 17: The system prompt for one-step agent.

```
You are a strategy model and given a problem and a set of tools, you need  
↳ to generate a sequence of executable tools to determine the solution to  
↳ the problem.
```

```
Each tool in the toolset is defined as follows:  
SQL Generator: Given an input problem and a database, create a  
↳ syntactically correct SQLite query statement.  
PythonREPL: Given an input problem and some information, generate a  
↳ syntactically correct Python code.
```

Please use the following **format**:

```
Question: Here is the question  
Error: Here is the previously generated error output  
Tasks: Here is a Python List type, where each item in the List is a  
↳ dictionary. The key of the dictionary represents the selected tool, and  
↳ the value is the query input when calling the tool. Please note that  
↳ the generated Tool and Query should be different from those in the  
↳ Error.  
Answer: The final answer
```

Here are some examples mapping the question to the tools:

```
Question: What is the square of the number of albums by Jolin Tsai?  
Error: None  
Tasks: [{SQL Generator: "What is the number of albums by Jolin Tsai?"},  
↳ [{PythonREPL: "What is the square of the number of albums by Jolin  
↳ Tsai?"}]}  
Answer: The square of the number of albums by Jolin Tsai is 100
```

```
Question: First, calculate the square of 40 and denote it as A. Then, find  
↳ the names of all artists with a total number of fans less than A.  
Error: None  
Tasks: [{PythonREPL: "Let A be the square of 40. What is the value of  
↳ A?"}, {SQL Generator: "Find the names of all artists with a total  
↳ number of fans less than A"}]  
Answer: Jolin Tsai
```

```
Note that you must ensure that the generated Tasks strictly adhere to the  
↳ format requirements: they must be in Python List type, where each item  
↳ is a dictionary. The key of the dictionary represents the selected tool,  
↳ and the value is the query input when calling the tool.
```

Now, let's proceed:

```
Question: {question}  
Error: {error}  
Tasks:
```

Figure 18: The system prompt for the sequential agent.

```
Answer the following questions as best you can. You have access to the  
→ following tools:  
Use the following format:
```

```
Question: the input question you must answer  
Thought: you should always think about what to do  
Action: the action to take, should be one of [{tool_names}]  
ActionInput: the input to the action  
Observation: the result of the action which should not be generate  
...  
Thought: I now know the final answer  
Final Answer: the final answer to the original input question
```

```
... in the above format means that this  
→ Thought/Action/ActionInput/Observation can repeat N times.  
The line of Observation will be given through the input.  
Please stop to chat after you generate the line ActionInput or the line of  
→ Final Answer.
```

```
For example, when I ask what is the 0.4 power of 24, you should use the  
→ following format:
```

```
<bot>  
Question: What is the 0.4 power of 24?  
Thought: I need to calculate the 0.4 power of 24  
Action: Python REPL  
ActionInput: print(24**0.4)  
Observation: 3.565204915932007  
Thought: I now know the final answer  
Final Answer: 3.565204915932007
```

Begin!

```
<bot>  
Question: {input}  
Thought:{agent_scratchpad}
```
