

S-Agent: an Agent Collaborative Framework Inspired by the Scientific Methodology

Anonymous ACL submission

Abstract

001 An increasing number of advancements have
002 been accomplished in agents empowered by
003 Large Language Models (LLM), particularly
004 in resolving simple dialogue tasks. However,
005 existing agents still face intractable robustness
006 issues for solving complex tasks, encounter-
007 ing the cascading hallucinations induced by
008 multi-step invocations of LLM. Certain recent
009 studies utilize multi-step reasoning, planning
010 strategies, and domain workflows to improve
011 the success rate of complex tasks, yet they ne-
012 glect the scientific methodology that encom-
013 passes the accumulated wisdom derived from
014 centuries of scientific inquiry. Drawing inspi-
015 ration from the scientific methodology, we pro-
016 pose the S-Agent - an agent collaborative frame-
017 work meticulously designed to actively exper-
018 iment and refine theories based on the analy-
019 sis of experimental results, thereby enhancing
020 the deductive capabilities of LLMs and comple-
021 menting their inductive and communicative
022 strengths. Additionally, we introduce an innova-
023 tive parallel planning methodology, wherein
024 agents with identical roles collaborate to simul-
025 taneously address the same inquiry. Extensive
026 experiments demonstrate the effectiveness and
027 efficiency of our approach. Notably, we achieve
028 a new state-of-the-art 33.3% pass@1 accuracy
029 on the LeetCodeHardGym coding benchmark
030 and a relatively good 96.3% pass@1 on Hu-
031 manEval with GPT-4.

032 1 Introduction

033 Recently, significant advancements have been
034 achieved in the realm of problem-solving through
035 agents founded upon Large Language Models
036 (LLM). Existing studies incorporated step by step
037 reasoning & planning strategies (Yao et al., 2023b;
038 Shinn et al., 2023; Yao et al., 2023a; Zhou et al.,
039 2023), and used tools to extend agents' capabili-
040 ties(Wu et al., 2023; Yang et al., 2023; Shen et al.,
041 2023). These studies have demonstrated their ca-
042 pability to tackle uncomplicated dialogue tasks.

043 However, current agents continue to confront insur-
044 mountable challenges in terms of resilience when
045 it comes to resolving complex tasks, as they en-
046 counter the cascading hallucinations brought about
047 by the multi-step invocations of LLM.

048 In order to improve the proficiency of agents
049 in resolving complex tasks, some recent studies
050 employ domain expertise to guide agents towards
051 adhering to standard operating procedures (SOP)
052 (Hong et al., 2023; Qian et al., 2023). These studies
053 can only serve specialized applications as they rely
054 on domain-specific procedural knowledge, still lack
055 of principled guidance.

056 To address the aforementioned issues, we draw
057 inspiration from the scientific method¹ that encom-
058 passes the wisdom accumulated through centuries
059 of scientific exploration and has been validated
060 across various disciplines. The magnificent mod-
061 ern science usually adheres to an iterative paradigm:
062 deriving ideas from observations and constructing
063 hypotheses. These hypotheses are then subjected
064 to experimentation, with the outcomes serving as
065 observations that either validate or question the hy-
066 potheses. At the heart of this paradigm lies the
067 notion of "falsifiability," introduced by Karl Pop-
068 per(Popper, 1959), which asserts that there must
069 exist experimental findings capable of disproving
070 the hypotheses. Guided by this paradigm, theory
071 and experiment can mutually inform and enhance
072 scientific understanding.

073 In this paper, we propose the S-Agent, a multi-
074 agent framework where agents partake in dialogues
075 and collaborations inspired by the scientific method.
076 This framework encompasses crucial processes in-
077 cluding idea generation, experiment conduction,
078 and the discussion of results. Our experiments

¹scientific method: A method of procedure that has char-
acterized natural science since the 17th century, consisting in
systematic observation, measurement, and experiment, and the
formulation, testing, and modification of hypotheses: criticism
is the backbone of the scientific method.

079 demonstrate that the system performs exceptionally
080 well in tackling challenging tasks, such as coding,
081 multi-hop QA .

082 To summarize, our key contributions are the fol-
083 lowing:

- 084 • To the best of our knowledge, we are pioneer-
085 ing the integration of formulating, experiment-
086 ing, and modifying of hypotheses within the
087 LLM agent system. These mechanisms repre-
088 sent the accumulated wisdom of centuries of
089 scientific research and are poised to enhance
090 the credibility and accuracy of LLM agents.
- 091 • We present S-Agent, a collaborative frame-
092 work for agents that integrates an automated
093 workflow planner and a parallel agent task
094 management unit. The framework offers
095 adaptable assistance for developing agents of
096 complex and high-reliability tasks, enabling
097 simultaneous operation at the agent level and
098 significantly enhancing efficiency.
- 099 • We conduct extensive experiments to demon-
100 strate the effectiveness and efficiency of
101 our approach. Notably, we achieve a new
102 state-of-the-art 33.3% pass@1 accuracy on the
103 leetcode-hard benchmark with GPT-4 and re-
104 lative good result 96.3% pass@1 accuracy on
105 the HumanEval coding benchmark with GPT-
106 4.

107 2 Related work

108 2.1 Think Like Human

109 In the early phase of the Language Model
110 (LLM) era, researchers began exploring LLMs
111 with the goal of achieving universal question-
112 answering capabilities, formulated as $answer =$
113 $LLM(question)$ (Devlin et al., 2018; Raffel et al.,
114 2020). Subsequently, GPT-3(Brown et al., 2020),
115 a pioneer in this domain, introduced the con-
116 cept of few-shot learning. This method reflects
117 the human ability to improve problem-solving
118 skills through exposure to a limited number
119 of examples, thereby shifting the paradigm to
120 $answer = LLM(demos, question)$. Further-
121 more, the Chain of Thought (COT) approach(Wei
122 et al., 2022) demonstrated that incorporating sim-
123 ple instructions such as 'please do it step by
124 step' significantly enhances performance. Con-
125 sequently, the paradigm evolved to $answer =$
126 $LLM(instructions, demos, question)$.

127 By crafting diverse instructions, the Tree of
128 Thoughts (ToT) framework(Yao et al., 2023a)
129 enables LLMs to generate multiple plans and
130 select the most appropriate one. The ReAct frame-
131 work(Yao et al., 2023b), which integrates reasoning
132 and action, modifies the paradigm to $result_{i+1} =$
133 $LLM(instruction, demos, result_i, question)$.
134 Reflexion(Shinn et al., 2023) advances this ap-
135 proach by supervising the entire decision-making
136 process and offering feedback on the complete
137 sequence of actions.

138 This paradigm closely mirrors the human pro-
139 cess of problem-solving: step-by-step refinement
140 of solutions based on past experiences and logical
141 analysis. However, we observe that when faced
142 with more complex issues, humans employ another
143 cognitive approach not yet utilized in LLM-based
144 agents: the scientific method, which includes con-
145 ducting experiments. Scientists design experiments
146 to test their theories, comparing experimental re-
147 sults with expectations to identify discrepancies
148 before actual implementation. Current LLM agents
149 lack this capability; they do not formulate expecta-
150 tions prior to taking actions.

151 2.2 Use Tool Like Human

152 While the Marxist philosophy posits that tool uti-
153 lization and creation are key distinctions between
154 humans and animals, it's the transformative impact
155 of tools on human evolution and dominance on
156 Earth that is truly noteworthy. This concept finds a
157 parallel in the realm of language models like Chat-
158 GPT. Despite their impressive performance and
159 global recognition, these models have inherent lim-
160 itations, including constrained calculation abilities,
161 restricted access to rare knowledge, and limited
162 proficiency in handling other modalities. Mirroring
163 the human approach to overcoming similar con-
164 straints, the strategic use of tools has emerged as
165 an effective solution in augmenting these models.

166 Pioneering efforts such as Visual ChatGPT(Wu
167 et al., 2023) and HuggingGPT(Shen et al., 2023)
168 laid the groundwork by integrating multi-modal
169 models as auxiliary tools, thereby expanding the
170 functionalities of LLMs beyond single-modal ca-
171 pabilities. This trajectory was further propelled
172 by MM-ReAct(Yang et al., 2023), which cleverly
173 incorporated a search engine and Microsoft API
174 services into the mix. This innovative approach
175 has been widely adopted, with OpenAI's ChatGPT
176 introducing plugin functions and Microsoft's New-
177 Bing exemplifying an LLM integrated with Bing

178 search capabilities.

179 In today’s landscape, the ability to configure
180 and customize tool sets has become a fundamental
181 feature in most agent assistant applications. In-
182 tegrating tools transforms the paradigm into $tool =$
183 $LLM(instructions, toolset, question), answer =$
184 $LLM(demos, tool, question)$. However, most
185 previous work has primarily focused on the
186 correctness of tool selection, without considering
187 how to enable LLMs to use complex tools. We aim
188 for our approach to empower LLMs to effectively
189 utilize tools with which they are not yet familiar.

190 2.3 Collaborate like Human

191 Researchers have ventured beyond exploring hu-
192 man cognition and tool usage, delving into the
193 intricacies of organizational dynamics. This ex-
194 ploration has been facilitated by employing multi-
195 agent systems to create artificial entities that sim-
196 ulate the workings of a company. The innovative
197 concept of role-playing was pioneered by Camel(Li
198 et al., 2023), utilizing the paradigm $answer =$
199 $LLM(..., personality, question)$.

200 This concept was expanded in subsequent
201 studies that modeled their operations explicitly
202 after a corporate structure. Notable examples
203 include MetaGPT(Hong et al., 2023) and Chat-
204 Dev(Qian et al., 2023), which emulate software
205 companies handling programming tasks. These
206 models assign roles like CEO, CTO, product
207 manager, programmer, and designer, system-
208 atically reflecting the organizational structure
209 of a real-world company. The adoption of a
210 Standard Operating Procedure (SOP), crafted
211 by professionals and fed into the system, guides
212 the collaborative process, creating a workflow
213 similar to that of a traditional software company.
214 These works shift the paradigm to $output_i =$
215 $LLM(personality_i, instruction_i, demos_i, input_i)$,
216 where each index represents a different role. By
217 following a sequence of roles, the answer is
218 generated through this process.

219 Building on this foundation, AgentVerse(Chen
220 et al., 2023b) and AutoAgents(Chen et al., 2023a)
221 introduced a job market system, simulating the re-
222 cruitment of expert agents for specific roles. This
223 approach also generate the SOP. By automating ev-
224 erything, the artificial company of agents achieves
225 a high level of task proficiency and autonomy.

226 While these models typically employ a sequen-
227 tial waterfall workflow, addressing each compo-
228 nent of a task linearly, recent research(Zhang et al.,

229 2023; Chen et al., 2023b; Du et al., 2023) has high-
230 lighted the benefits of cooperative approaches, such
231 as debates, in enhancing performance. Our work
232 introduces an innovative parallel planning method-
233 ology, where agents with identical roles collaborate
234 to simultaneously address the same question. This
235 parallel approach has been instrumental in boosting
236 performance, demonstrating the efficacy of multi-
237 agent cooperation in complex problem-solving.

238 3 Methodology

239 To empower the S-Agent system with the scientific
240 method, we have structured the process into three
241 distinct parts: idea generation, experiment conduc-
242 tion, and panel discussion. Each part is facilitated
243 by purpose-built agents, as illustrated in figure 1.

244 3.1 Idea Generation

245 The Idea generation part serves as the initial stage
246 for proposing ideas and formulating hypotheses
247 during the execution process of the agent system.
248 Throughout the experiment iterations, ideas un-
249 dergo refinement and new ones emerge. In this
250 phase, we specifically design LLM-powered agents,
251 termed **idea generation agents**, to emulate the role
252 of scientists in generating structured ideas easily
253 verifiable by subsequent experiments. Meanwhile,
254 these agents can analyze feedback from discussion
255 parts and revise their ideas accordingly. In pro-
256 gramming scenarios, these agents initially receive
257 a coding question as input and directly generate
258 Python solutions. In the next few rounds, the agent
259 shall take the discussion feedback as input and gen-
260 erate new solutions if the solutions generated by the
261 previous ideas do not pass the targeted experiments.
262 As shown in Figure 1, the idea generation agents
263 receive the input and produce ideas. A detailed
264 task adaption will be presented in the experiments
265 section.

266 3.2 Experiments Conduction

267 The experiment conduction phase consists of de-
268 signing appropriate experiments and carrying out
269 the designed experiments. These steps are often
270 considered the most crucial in the scientific method.
271 A well-designed experiment can determine not only
272 the validity of a hypothesis but also highlight its
273 advantages over previous theories when the results
274 support the hypothesis. Conversely, when the re-
275 sults reject the hypothesis, the experiment can pin-
276 point the specific issues that remain unresolved.
277 This focused feedback can make the iteration of

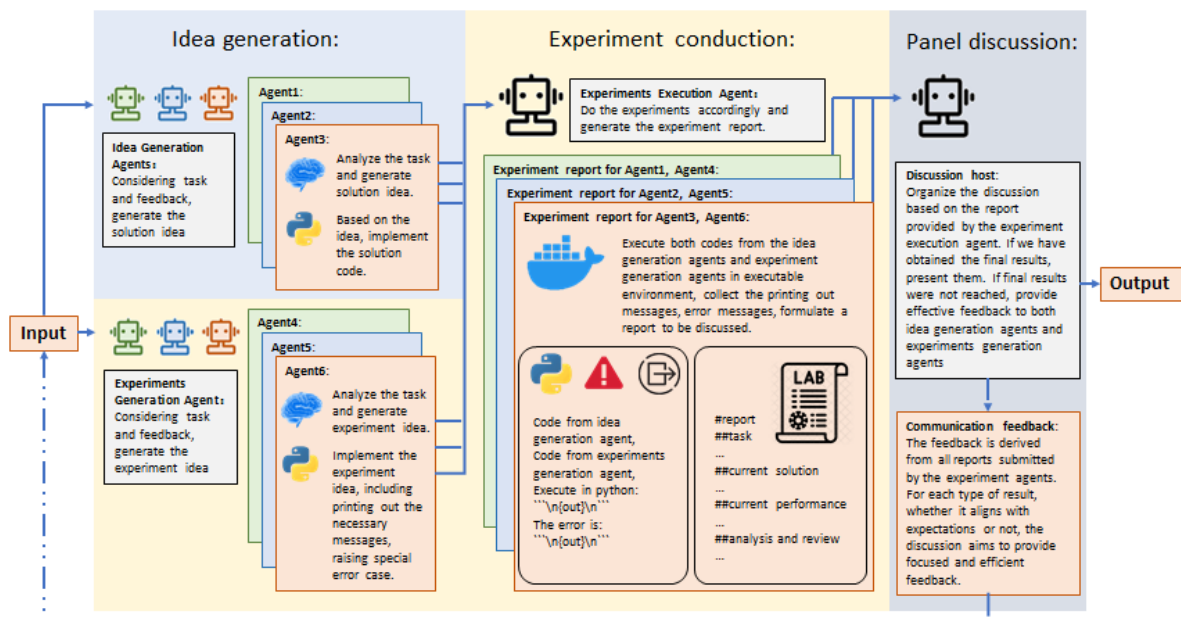


Figure 1: Illustration of S-Agent. The idea generation phase showcases agents responsible for generating ideas, encompassing both solution concepts and corresponding code. The experiment conduction phase begins with the generation of experiment ideas by experiment generation agents, followed by the execution of experiments through experiment execution agents. Experiment reports contained thorough analysis of results and the collection of feedback from the environment are also generated in this phase. In the panel discussion session, a discussion host synthesizes experiment reports from preceding agents, facilitating the generation of valuable feedback. This feedback is systematically organized to serve as input for the subsequent iteration, fostering continuous idea and experiment setting improvement.

278 hypotheses much faster, also help avoid deduction
 279 analysis from wrong start. In our system, applying
 280 this hypotheses experiment alignment phase properly
 281 can avoid machine hallucination effectively.
 282 To equip this mechanism in our system, we designed
 283 **experiment generation agents** to do experiments
 284 design and **experiment execution agents**
 285 to test out hypotheses and produce detailed exper-
 286 iment reports after receiving the results. As illus-
 287 trated in Figure 1, the experiment generation agents
 288 formulated the targeted experiment plan based on
 289 proposed idea and execution reports, and the exper-
 290 iment execution agents execute codes with these
 291 test cases in execution environment.

292 3.3 Panel Discussion

293 Observing the workflow of a scientist reveals that
 294 panel discussion, often through paper publications
 295 and sharing, is a vital component of the modern
 296 academic community. Therefore, when an exper-
 297 iment validates an idea, we facilitate discussion
 298 among different agents. This process allows them
 299 to review others' work, ensuring internal logic con-
 300 sistency and providing references and inspiration

to other agents. The discussion results in either
 feedback on analyzing experiment results or final
 answers if the designed experiments successfully
 verify the proposed ideas. To manage this pro-
 cess effectively, we have designed a specialized
 agent called a **discussion host**, responsible for ag-
 gregating information, assessing previous results,
 and overseeing the overall status of ongoing discus-
 sions.

3.4 Supporting Component

To simplify the deployment of the entire system and
 reduce execution time, we also include two support-
 ing components named the planner and agent task
 management.

Automated Workflow Planner The automated
 workflow planner is a specially designed agent re-
 sponsible for generating sequences of execution
 flows for agents and managing their interdependen-
 cies. This critical component strategically plans
 and organizes the order in which agents operate,
 ensuring a coherent and efficient workflow. To for-
 mulate this plan, only the input-output information

and the agents’ goals are required. As shown in Figure 1, the edges of the directed workflow graph are generated by the planner automatically. This approach draws inspiration from the methodology introduced in LLMCompiler(Kim et al., 2023). The detailed prompt design and sample illustration are given in the appendix.

Agent Task Management Unit The agent task management unit, drawing inspiration from the instruction fetching mechanism in modern computer architecture, plays a crucial role in determining the optimal execution flow of agents based on the intermediate representation generated by the Planner. Employing a greedy policy, this unit swiftly adds agents to the task list as soon as they become ready for parallel calling. The implementation of this agent task management unit involves a straightforward fetching and queue mechanism, foregoing the need for a dedicated agent system.

4 Experiment

We assess the effectiveness of our framework in tackling complex problems using the coding benchmark LeetcodeHardGym(Shinn et al., 2023), alongside evaluations through HumanEval (Chen et al., 2021) and EvalPlus(Liu et al., 2023). Furthermore, we conduct a focused analysis to evaluate the efficacy of our approach in reducing hallucinations on the multi-hop QA benchmark HotpotQA dataset (Yang et al., 2018).

4.1 Coding

HumanEval(Chen et al., 2021)(distributed under the MIT license) is a benchmark for code synthesis, which consists of 164 programming problems with several test cases each. The problems in this dataset are designed to test the ability of LLMs to generate functionally correct codes, which means the generated codes can not only execute successfully but also pass the provided test cases, instead of being linguistically similar to the canonical solution. EvalPlus is a benchmark that aims to improve the quality and quantity of test cases for the existing HumanEval benchmark. EvalPlus contains new test cases that can catch more errors and bugs in the LLM-synthesize code. The LeetCodeHardGym(Shinn et al., 2023) consists of 40 open LeetCode hard problems. It was introduced in Reflexion, where the benchmark utilizes LeetCode’s API and the traditional RL package gym (Brockman et al., 2016) to construct this dataset in the

humaneval format, requiring no additional configuration modifications.

Our performance on these benchmarks is noteworthy, achieving a Pass@1 (Pass@1 is the probability that a model generates at least one correct solution out of one attempt) of 33.3% on LeetCodeHardGym, establishing a new SOTA. Additionally, we achieve strong performance with 96.3% on HumanEval and 86.6% on EvalPlus,

4.1.1 Implementation Details

As discussed earlier, the scientific method involves three phases: idea generation, experiment conduction, and panel discussion. In this section, we detail our approach to applying this methodology in the adaptation of the HumanEval dataset. We elaborate on our process for prompt design and provide a detailed example in the appendix for clarity.

In the idea generation phase, we task the idea generation agents with the dual roles of analysis and coding. Initially, these agents analyze and formulate a comprehensive strategy to tackle the current task. Subsequently, during the coding phase, they annotate each step of their proposed solution with explanatory comments. This process can be represented as $analysis, modified_code = LLM(task, previous_code, feedback)$. This approach mirrors the scientific method where $analysis, modified_theory = Scientist(phenomenon, previous_theory, experiment_result)$ is analogous. The general structure of the prompt format is detailed in Figure 2, which also specifies the output format.

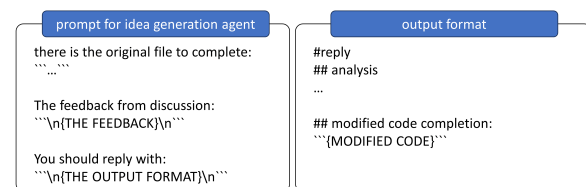


Figure 2: Prompt for idea generation agents and the output template. FEEDBACK is all the previous experiments reports.

During the experiment conduction stage, specialized agents are designated for experiment generation and execution. Experiment generation agents are tasked with creating specific experiments (in the context of coding, experiments refer to test cases), determining the experimental inputs, defining expected outputs, specifying messages for various types of output during analysis, and providing implementation code formatted

as 'assert f(input)==output, description of the experiment'. Detailed instructions for this process are illustrated in Figure 4, where $test_script = LLM(task, current_code)$. This parallels the scientific method, where $experiment_plan = experiment_scientist(phenomenon, current_theory)$.

The experiment execution agent is tasked with running Python code generated from idea generation within a functional environment and generating the experiment report. The output format of this agent is outlined in Figure 5, with detailed specifications of the report format provided in Figure 10 in the appendix. The agent interprets standard output and error messages to provide precise feedback, facilitating the refinement of theories or experiment designs.

Following the generation of experiment reports during the experiment conduction phase, the discussion host assumes responsibility for aggregating all reports and assessing the resolution of the issue. If the issue persists, each idea generator conducts a thorough analysis of their respective reports to refine their proposed ideas embedded within the generated code. Moreover, each participant involved in experiment generation reviews the reports to make adjustments to their experiments. This process exemplifies the effective panel discussion mechanism established by our framework.

In each test case, we initiate by feeding input to our system. Subsequently, each interaction with received feedback is considered a round. We set a maximum of 10 rounds per test case; exceeding this limit leads to the system being deemed unsuccessful. Conversely, upon achieving success, we derive the final solution from the resulting code. Upon completing all 164 cases, we run the official script to compute a Pass@1 score. We chose the number 10 to balance the system's performance, as increasing this limit might lead to chance improvements without substantial benefit. Comparable frameworks, such as Reflexion(Shinn et al., 2023), use a parameter of 6 for similar reasons. Additionally, our analysis shows that among the 158 successful cases, only 4 required more than 2 rounds to solve. Hence, reducing this number would not notably affect overall performance.

4.1.2 Result and Analysis

We evaluated S-Agent system using both GPT-4 and GPT-3.5. By utilizing GPT-4, our system successfully completed 158 out of 164 tasks (pass@1

= 96.3%), achieving state-of-the-art performance. Furthermore, our system leveraged GPT-3.5 to pass 137 out of 164 tasks (pass@1 = 84.1%), surpassing the performance of all our known works based on the same LLM. The results are summarized in Table 1. We have also attained a SOTA result on the LeetCodeHard benchmark, achieving a Pass@1 accuracy of 33%. Previously, the highest accuracy was 15%, achieved by Reflexion(Shinn et al., 2023).

Recent methodologies in coding tasks aim to enhance performance through several strategies: leveraging the Python execution environment, incorporating reflection mechanisms, and facilitating multi-agent discourse. The Python execution environment allows systems to execute code and verify its validity. Reflection enhances system performance by using verbal reinforcement cues generated by LLMs based on past experiences. Different systems implement reflection in varied ways: MetaGPT(Hong et al., 2023) employs a specialized agent for reflection, language-agent-tree-search (LATS)(Zhou et al., 2023) utilizes reflection in a trajectory format, and Reflexion(Shinn et al., 2023) focuses on enabling the LLM itself to review and learn from feedback. In coding tasks, feedback includes outputs such as standard output (stdout) and error messages generated during code execution. Multi-agent discourse, studied extensively in this context, involves aggregating ideas from different personas to generate more comprehensive and accurate answers(Du et al., 2023; Chen et al., 2023b).

Prior methodologies typically receive feedback passively from the Python Interpreter. In contrast, our approach uniquely empowers the agent to actively influence the feedback process, as detailed in Figure 3. In our framework, outlined in the implementation section, the experiment generation agents configure experiments to produce intermediate results. This approach introduces additional lines such as printed statements that provide detailed insights into any encountered errors. These lines not only utilize feedback but also enrich it through agent-driven efforts, thereby enhancing the feedback with more valuable information. This enrichment guides the system towards effectively addressing complex problems. The enhanced feedback from experimental results and inter-agent panel discussions has propelled us to achieve state-of-the-art performance.

	GPT4	ANPL	MetaGPT	AgentVerse	Reflexion	LATS	S-Agent
code execution	×	✓	✓	✓	✓	✓	✓
reflection	×	×	✓	✓	✓	✓	✓
panel discussion	×	×	×	✓	×	×	✓
experiment design	×	×	×	×	×	×	✓
Pass@1 (GPT3.5-base)	-	76.2%	-	75%	-	83.8%	84%
(GPT4-base)	67%	86.6%	87.7%	89%	91%	94.4%	96.3%

Table 1: Pass@1 result of related works on programming. We refer code execution as the use of python execution environment, reflection as the use of LLM-generated feedback, panel discussion as agents sharing and discuss each others work, experiment design as actively design experiment according to proposed idea and previous feedback. Our S-Agent system stands out from previous works and achieved best pass@1 score, mainly because of involving actively designing experiments.

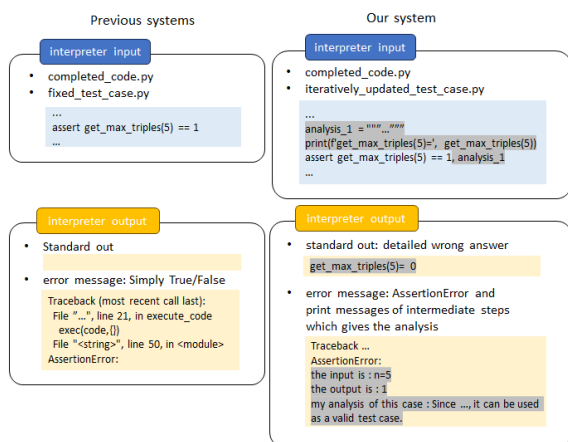


Figure 3: This figure demonstrates the differences in input and output between our system and the previous systems. Because of incorporating specially designed experiments during the python execution step, instead of just providing a binary assessment of code functionality, our system also generate intermediate results and other results based on the experiment requirement, as highlighted in the image.

4.1.3 Ablation Study

We conduct experiments on the HumanEval and EvalPlus datasets to investigate the effectiveness of mechanisms within the S-Agent framework.

Table 2 shows the results of the experimental results of three S-Agent variants including the original framework (Original), S-Agent without Panel discussion mechanism (w/o Panel), S-Agent without Panel discussion and Experiment conduction mechanism (w/o Exp&Panel), S-Agent without the Idea generation, Experiment conduction, and Panel discussion mechanisms (w/o Idea&Exp&Panel). The findings indicate that all three fundamental mechanisms play a beneficial role in enhancing the precision of the S-Agent framework.

We also investigate how the number of Idea gen-

S-Agent Variants	HumanEval	EvalPlus
Original	96.3%	89.0%
w/o Panel	92.1%	86.6%
w/o Exp&Panel	84.8%	79.9%
w/o Idea&Exp&Panel	68.3%	65.2%

Table 2: Pass@1 on HumanEval and EvalPlus using GPT-4 under different S-Agent Variants. The experiment employs a single agent for idea generation, another for experiment conduction, and a third for panel discussion.

eration and Experiment conduction agents impact the results of the S-Agent framework. As depicted in Table 3, with an increase in the number of agents, there is a clear improvement in the average pass@1 value, accompanied by a reduction in the standard deviation. This indicates that, solely from a performance perspective, increased discussion can substantially enhance both the accuracy and stability of the entire system. This conclusion also aligned with (Du et al., 2023; Yang et al., 2018). This also aligned with our result in GPT-4 based experiments, the pass@1 increase from 92.1% to 96.3% when increasing the number of agents to 3.

Number of Agents	HumanEval Pass@1 (%)	EvalPlus Pass@1 (%)
1	80.5 ± 1.8	72.0 ± 1.8
3	83.5 ± 1.37	74.4 ± 1.37
5	84.1 ± 0.56	75.6 ± 0.56

Table 3: Pass@1 performance on HumanEval and EvalPlus using GPT-3.5 under different numbers of Idea generation and Experiment conduction agents.

4.2 HotpotQA

The HotpotQA(Yang et al., 2018) dataset, is a crucial benchmark in Natural Language Processing (NLP) and Question Answering (QA). It is specifically designed for models that handle complex, multi-hop question-answering tasks, requiring synthesis of information across multiple text sources from wiki. The HotpotQA dataset is distributed under the CC BY-SA 4.0 license. The code is distributed under the Apache 2.0 license.

We randomly sampled a subset of 50 QA pairs from the original dataset and conducted tests without panel discussion settings. Compared to ReAct(Yao et al., 2023b), our approach improved results by 7.4%. The S-Agent outperforms ReAct(Yao et al., 2023b) by providing more precise and effective feedback.

In our experiment demonstrating how LLM’s approach can reduce hallucinations, we examined a query asking "Which show was hosted by Jessica Drake’s former spouse?". Both our approach and previous methods use the current result to generate subsequent queries based on Wikipedia. The key difference lies in how queries are generated: previous methods employ $next_query = LLM(task, current_knowledge, review_of_lastquery)$, whereas our approach utilizes $next_query = LLM(task, current_knowledge, expanded_of_last_query, last_query_result)$

In this scenario, both systems first query information about Jessica Drake’s ex-husbands, identifying Brad Armstrong and Evan Stone. The React system, lacking explicit expectation management, falls into a loop when querying Brad Armstrong. It assumes the direction is correct solely based on his association with the movie industry, perpetuating the query without finding the answer. In contrast, our framework treats each query as an experiment with predefined expectations. By comparing the actual result with these expectations, LLM identifies incorrect paths and intervenes with feedback such as "The provided text does not mention any show hosted by Brad Armstrong". This feedback redirects attention to Evan Stone, leading to successful resolution of the task.

The precise and efficient feedback mechanism in our framework plays a crucial role in preventing LLM from persisting in incorrect directions, thereby mitigating the risk of endlessly generating inaccurate information. For detailed prompts,

please refer to the appendix.

5 Conclusion

In this paper, we introduce the S-Agent, an innovative multi-agent framework in which agents engage in dialogues and collaborations inspired by the scientific methodology. The framework incorporates the essential processes of hypothesis development, experimentation, and refinement. These processes embody the collective knowledge accumulated over centuries of scientific inquiry and are poised to enhance the credibility and precision of LLM agents. The S-Agent integrates an automated workflow planner and a parallel agent task management unit, providing flexible support for developing agents for complex and high-stakes tasks, facilitating concurrent operation at the agent level. Extensive experiments confirm the efficacy and efficiency of our methodology. Notably, the S-Agent achieves a new state-of-the-art 96.3% pass@1 accuracy on the HumanEval coding benchmark with GPT-4.

6 Limitations and Future Work

The primary objective of this paper is to elucidate the concept of scientific methodology. While our general framework may not exhibit the same level of sophistication as pioneering works like AutoAgents(Chen et al., 2023a) which autonomously generate requisite agents, our framework still requires the manual implementation of specific agents tailored for specialized tasks. In future developments, a key direction of research involves exploring methods to automatically generate agents with finely crafted prompts, presenting an important avenue for further exploration.

Currently, following the generation of the Directed Acyclic Graph (DAG) plan, the plan remains static. However, it is imperative to establish dynamic refinement for this plan. Recent advancements, exemplified by works such as ReAct(Yao et al., 2023b), BabyAGI, and XAgent, have endeavored to enhance plans based on feedback received at each step. While these approaches typically involve linearly designed steps, there is a research gap in developing methods to dynamically refine a DAG-formatted plan with the capability for parallel execution.

References

- 642 Greg Brockman, Vicki Cheung, Ludwig Pettersson,
643 Jonas Schneider, John Schulman, Jie Tang, and Woj-
644 ciech Zaremba. 2016. [Openai gym](#).
- 645 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
646 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
647 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
648 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
649 Gretchen Krueger, Tom Henighan, Rewon Child,
650 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens
651 Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-
652 teusz Litwin, Scott Gray, Benjamin Chess, Jack
653 Clark, Christopher Berner, Sam McCandlish, Alec
654 Radford, Ilya Sutskever, and Dario Amodei. 2020.
655 [Language models are few-shot learners](#). In *Ad-
656 vances in Neural Information Processing Systems*,
657 volume 33, pages 1877–1901. Curran Associates,
658 Inc.
- 659 Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Sesay
660 Jaward, Karlsson Börje, Jie Fu, and Yemin Shi. 2023a.
661 Autoagents: The automatic agents generation frame-
662 work. *arXiv preprint*.
- 663 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming
664 Yuan, Henrique Ponde de Oliveira Pinto, Jared Kap-
665 lan, Harri Edwards, Yuri Burda, Nicholas Joseph,
666 Greg Brockman, Alex Ray, Raul Puri, Gretchen
667 Krueger, Michael Petrov, Heidy Khlaaf, Girish Sas-
668 try, Pamela Mishkin, Brooke Chan, Scott Gray,
669 Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz
670 Kaiser, Mohammad Bavarian, Clemens Winter,
671 Philippe Tillet, Felipe Petroski Such, Dave Cum-
672 mings, Matthias Plappert, Fotios Chantzis, Eliza-
673 beth Barnes, Ariel Herbert-Voss, William Hebgen
674 Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie
675 Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain,
676 William Saunders, Christopher Hesse, Andrew N.
677 Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan
678 Morikawa, Alec Radford, Matthew Knight, Miles
679 Brundage, Mira Murati, Katie Mayer, Peter Welinder,
680 Bob McGrew, Dario Amodei, Sam McCandlish, Ilya
681 Sutskever, and Wojciech Zaremba. 2021. [Evaluating
682 large language models trained on code](#).
- 683 Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang,
684 Chenfei Yuan, Chen Qian, Chi-Min Chan, Yujia
685 Qin, Yaxi Lu, Ruobing Xie, et al. 2023b. [Agent-
686 verse: Facilitating multi-agent collaboration and ex-
687 ploring emergent behaviors in agents](#). *arXiv preprint
688 arXiv:2308.10848*.
- 689 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
690 Kristina Toutanova. 2018. Bert: Pre-training of deep
691 bidirectional transformers for language understand-
692 ing. *arXiv preprint arXiv:1810.04805*.
- 693 Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenen-
694 baum, and Igor Mordatch. 2023. [Improving factual-
695 ity and reasoning in language models through multia-
696 gent debate](#). *arXiv preprint arXiv:2305.14325*.
- 697 Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu
698 Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin Wang,
Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang
Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu,
and Jürgen Schmidhuber. 2023. [Metagpt: Meta pro-
gramming for a multi-agent collaborative framework](#).
- Sehoon Kim, Suhong Moon, Ryan Tabrizi, Nicholas
Lee, Michael Mahoney, Kurt Keutzer, and Amir Gho-
lami. 2023. [An llm compiler for parallel function
calling](#). *arXiv*.
- Guohao Li, Hasan Abed Al Kader Hammoud, Hani
Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023.
[Camel: Communicative agents for "mind" explo-
ration of large language model society](#). In *Thirty-
seventh Conference on Neural Information Process-
ing Systems*.
- Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Ling-
ming Zhang. 2023. [Is your code generated by chat-
GPT really correct? rigorous evaluation of large lan-
guage models for code generation](#). In *Thirty-seventh
Conference on Neural Information Processing Sys-
tems*.
- Karl Popper. 1959. *The Logic of Scientific Discovery*,
2002 pbk; 2005 ebook edition. Routledge.
- Chen Qian, Xin Cong, Wei Liu, Cheng Yang, Weize
Chen, Yusheng Su, Yufan Dang, Jiahao Li, Juyuan
Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023.
[Communicative agents for software development](#).
- Colin Raffel, Noam Shazeer, Adam Roberts, Kather-
ine Lee, Sharan Narang, Michael Matena, Yanqi
Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the
limits of transfer learning with a unified text-to-text
transformer](#). *Journal of Machine Learning Research*,
21(140):1–67.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li,
Weiming Lu, and Yueting Zhuang. 2023. [Hugging-
gpt: Solving ai tasks with chatgpt and its friends in
huggingface](#). In *Advances in Neural Information
Processing Systems*.
- Noah Shinn, Federico Cassano, Edward Berman, Ash-
win Gopinath, Karthik Narasimhan, and Shunyu Yao.
2023. [Reflexion: Language agents with verbal rein-
forcement learning](#).
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le,
and Denny Zhou. 2022. [Chain of thought prompt-
ing elicits reasoning in large language models](#). In
Advances in Neural Information Processing Systems.
- Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong
Wang, Zecheng Tang, and Nan Duan. 2023. [Visual
chatgpt: Talking, drawing and editing with visual
foundation models](#).
- Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin
Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu,
Ce Liu, Michael Zeng, and Lijuan Wang. 2023. [Mm-
react: Prompting chatgpt for multimodal reasoning
and action](#).

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023a. [Tree of Thoughts: Deliberate problem solving with large language models](#).

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023b. ReAct: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*.

Jintian Zhang, Xin Xu, and Shumin Deng. 2023. [Exploring collaboration mechanisms for llm agents: A social psychology view](#).

Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. 2023. [Language agent tree search unifies reasoning acting and planning in language models](#).

A Additional Details

A.1 Automated Workflow Planner

Given a user query, the planner will create a plan to solve the user query with the most parallelizability. The prompt of the planner agent contains the profiles of every agent and the tools it equips. Each task in the plan must have a unique ID, which is strictly increasing. Inputs for each task can either be constants or outputs from preceding actions. If the outputs from preceding actions are needed, we use the format \$id to denote the ID of the previous agent task whose output will be replaced with the input. Upon the completion of all agent tasks, we always call join as the last task in the plan to collect all the previous task outputs and formulate a final output. We use the HumanEval coding question *get_max_triplets* to demonstrate the functionality of our LLM planner agents. Figure 6 shows the plan that the planner creates after receiving the coding question as the user query. In this case, user queries are fed into three Coders and the Tester. Afterward, the output from the Coders should be considered as the input to the Experimenters with the output from the Tester respectively. Then the gathered outputs from these Experiments should be the input of the Discussion Hoster. In the end, we collect all the results and finish this plan.

A.2 Agent Task Management Unit

Figure 7 shows how the management unit handles the agent tasks. Agents are equipped with the tools that the user provides and tasks are delegated to the associate tool. The management unit synchronously listens to the task queue and schedules tasks as they arrive in the queue based on their dependencies. More specially, in this case, the three coders and the Tester agents execute in parallel at the same time due to empty dependencies. While the Experiment agents cannot execute in parallel until the completion of all Coder agents and Tester Agents. Meanwhile, we shall replace dependency placeholders, i.e. \$i, in the args of the agent task with the actual output.

B Additional Results

B.1 Case study for HumanEval

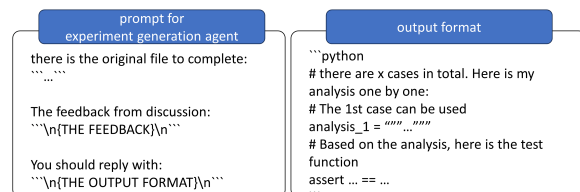


Figure 4: Prompt for experiments generation agent and the output template. FEEDBACK is all the previous experiments reports.

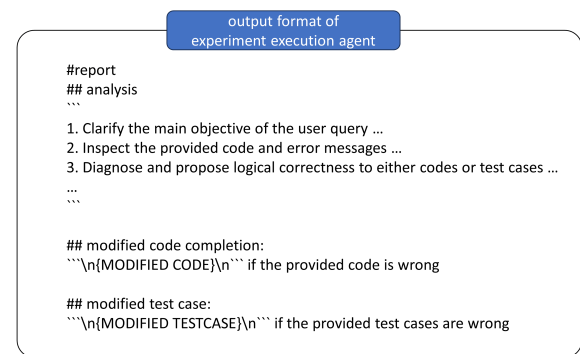


Figure 5: Output format of the experiment execution agent, where detailed analysis and modification of either codes or test cases are produced.

To get a closer look at a specific task, we chose No.147 as a demo because this task is only solved by our system, and it cost 5 rounds of modification to get the final answer. For the task and idea generation agent, refer figure 8, for the experiment generation agent refer figure 9, for the experiment execution agent and the exact augment feedback from python, refer figure 10

1. Coder_1(def get_max_triplets(n): """ You are given a positive integer n. You have to create an integer array a of length n. For each i (1 ≤ i ≤ n), the value of a[i] = i * i - i + 1. Return the number of triples (a[i], a[j], a[k]) of a where i < j < k, and a[i] + a[j] + a[k] is a multiple of 3. Example : Input: n = 5 Output: 1 Explanation: a = [1, 3, 7, 13, 21] The only valid triple is (1, 7, 13). """)
2. Coder_2(def get_max_triplets(n): """ You are given a positive integer n. You have to create an integer array a of length n. For each i (1 ≤ i ≤ n), the value of a[i] = i * i - i + 1. Return the number of triples (a[i], a[j], a[k]) of a where i < j < k, and a[i] + a[j] + a[k] is a multiple of 3. Example : Input: n = 5 Output: 1 Explanation: a = [1, 3, 7, 13, 21] The only valid triple is (1, 7, 13). """)
3. Coder_3(def get_max_triplets(n): """ You are given a positive integer n. You have to create an integer array a of length n. For each i (1 ≤ i ≤ n), the value of a[i] = i * i - i + 1. Return the number of triples (a[i], a[j], a[k]) of a where i < j < k, and a[i] + a[j] + a[k] is a multiple of 3. Example : Input: n = 5 Output: 1 Explanation: a = [1, 3, 7, 13, 21] The only valid triple is (1, 7, 13). """)
4. Tester(def get_max_triplets(n): """ You are given a positive integer n. You have to create an integer array a of length n. For each i (1 ≤ i ≤ n), the value of a[i] = i * i - i + 1. Return the number of triples (a[i], a[j], a[k]) of a where i < j < k, and a[i] + a[j] + a[k] is a multiple of 3. Example : Input: n = 5 Output: 1 Explanation: a = [1, 3, 7, 13, 21] The only valid triple is (1, 7, 13). """)
5. Experimenter(\$1+++4)
6. Experimenter(\$2+++4)
7. Experimenter(\$3+++4)
8. DiscussionHoster(\$5+++6+++7)
9. join()

Figure 6: Given a coding task: "def *get_max_triplets*: You are given a positive integer *n*. You have to create an integer array *a* of length *n*. For each *i* ($1 \leq i \leq n$), the value of $a[i] = i \times i - i + 1$. Return the number of triples ($a[i], a[j], a[k]$) of *a* where $i < j < k$, and $a[i] + a[j] + a[k]$ is a multiple of 3. Example: Input: $n = 5$, Output: 1, Explanation: $a = [1, 3, 7, 13, 21]$, The only valid triple is (1, 7, 13)", the planner agent shall automatically make the detailed plan of the agents' workflow and their dependencies.

Based on such augmented feedback, these agents can get more specific reviews about current code, just like scientists can draw more specific conclusions based on specific experiment results. Here they find out that the problem didn't handle reminder 1 and 2 properly.

Based on this feedback, they will repeat the work procedure. By observing the correct final answer, we would find that the issue at the first draft is that it did not consider three 1s or three 2s can also lead to reminding 0, just like three 0s.

B.2 HotpotQA

B.2.1 dataset

The HotpotQA (Yang et al., 2018) dataset stands as a pivotal benchmark within the realm of Natural Language Processing (NLP) and Question Answering (QA). This dataset serves a distinctive purpose, tailored for assessing models tasked with intricate, multi-hop question-answering assignments that necessitate the synthesis of information from various textual sources. Consequently, the inclusion of HotpotQA in evaluations requires models to showcase sophisticated reasoning and comprehension abilities.

HotpotQA offers two evaluation settings, Full-

wiki and Distractor. In the Fullwiki Setting, the dataset provides only questions, and users must retrieve related information from the entire Wikipedia dataset using their Information Retrieval (IR) system. The effectiveness of the search strategy in this setting is crucial, as the content found can significantly influence the results. Users then use their models to answer the question based on this information. In the Distractor Setting, the dataset provides questions along with context from the Wikipedia dataset, which includes both related and unrelated paragraphs. In this case, the user's model must be able to sift through the shuffled context to find the relevant information and answer the question correctly. In both settings, models are tasked with predicting the answer and identifying the supporting paragraphs in the context. When evaluating performance, we use Exact Match (EM), which measures whether the model's answer precisely matches the ground truth answer.

B.2.2 implementation

In this particular dataset, we empower idea generation agents to create the entire search plan. This includes formulating a set of queries for interaction with Wikipedia and specifying the desired

1. Task(idx=1, name='Coder_1', tool=<function writecode_1 at 0x11c11b370>, args=('\\ndef get_max_triplets(n):\\n """\\n You are given a positive integer n. You have to ..."""\n'), dependencies=[], stringify_rule=<function <lambda> at 0x11c11b2e0>, thought="", observation=None, is_join=False)
2. Task(idx=2, name='Coder_2', tool=<function writecode_2 at 0x11c11b520>, args=('\\ndef get_max_triplets(n):\\n """\\n You are given a positive integer n. You have to ..."""\n'), dependencies=[], stringify_rule=<function <lambda> at 0x11c11b6d0>, thought="", observation=None, is_join=False)
3. Task(idx=3, name='Coder_3', tool=<function writecode_3 at 0x11c11b760>, args=('\\ndef get_max_triplets(n):\\n """\\n You are given a positive integer n. You have to ..."""\n'), dependencies=[], stringify_rule=<function <lambda> at 0x11c11b9a0>, thought="", observation=None, is_join=False)
4. Task(idx=4, name='Tester', tool=<function create_unit_tests at 0x11c11bd00>, args=('\\ndef get_max_triplets(n):\\n """\\n You are given a positive integer n. You have to ..."""\n'), dependencies=[], stringify_rule=<function <lambda> at 0x11c13c0d0>, thought="", observation=None, is_join=False)
5. Task(idx=5, name='Experimenter', tool=<function experiment at 0x11c13fe20>, args=('\$1+++\$4',), dependencies=[1, 4], stringify_rule=<function <lambda> at 0x11c13ff40>, thought="", observation=None, is_join=False)
6. Task(idx=6, name='Experimenter', tool=<function experiment at 0x11c13fe20>, args=('\$2+++\$4',), dependencies=[2, 4], stringify_rule=<function <lambda> at 0x11c13ff40>, thought="", observation=None, is_join=False)
7. Task(idx=7, name='Experimenter', tool=<function experiment at 0x11c13fe20>, args=('\$3+++\$4',), dependencies=[3, 4], stringify_rule=<function <lambda> at 0x11c13ff40>, thought="", observation=None, is_join=False)
8. Task(idx=8, name='DiscussionHoster', tool=<function discuss at 0x11c13c1f0>, args=('\$5+++\$6+++\$7',), dependencies=[5, 6, 7], stringify_rule=<function <lambda> at 0x11c13c4c0>, thought="", observation=None, is_join=False)
9. Task(idx=9, name='join', tool=<function instantiate_task.<locals>.<lambda> at 0x11c35add0>, args=('<END_OF_PLAN>',), dependencies=[1, 2, 3, 4, 5, 6, 7, 8], stringify_rule=None, thought=("",), observation=None, is_join=True)

Figure 7: Here is an example of how the agent task management unit works for the coding task: `def get_max_triplets`, where tasks with no dependencies execute parallelly.

878 knowledge to be retrieved through these queries.
 879 Simultaneously generating pairs of queries and tar-
 880 geted knowledge streamlines the experimental pro-
 881 cess, especially when assessing the efficiency of
 882 the queries. This approach aligns with the scientific
 883 method, emphasizing the importance of providing
 884 falsifiable ideas.

885 In this context, the experiment is to evaluate
 886 whether the query can retrieve the targeted knowl-
 887 edge, therefore, a rule-based experiment generation
 888 agent is employed.

889 The experiment execution agent in this scenario
 890 is more intricate due to the complexity of com-
 891 municating with the Wikipedia API. During the
 892 implementation of this specific experiment execu-
 893 tion agent, we let llm automatically choose from
 894 several actions. If a Wikipedia page is success-
 895 fully retrieved by a query, the agent systematically
 896 analyzes the information, starting from the first
 897 paragraph and summary information box, followed
 898 by the table of contents (TOC) of the Wikipedia
 899 page. If the answer can be directly generated from
 900 this initial information, it is returned promptly. If
 901 no answer is found, but certain paragraphs poten-
 902 tially contain details relevant to answering the ques-
 903 tion, the agent continues to read these paragraphs.
 904 The answer is returned if discovered in this pro-
 905 cess; otherwise, a summary and analysis detailing
 906 the reasons for the failure to find an answer are

provided.

In situations where a Wikipedia page cannot be
 located with the given query, the agent generates a
 summary of similar pages returned by the API.

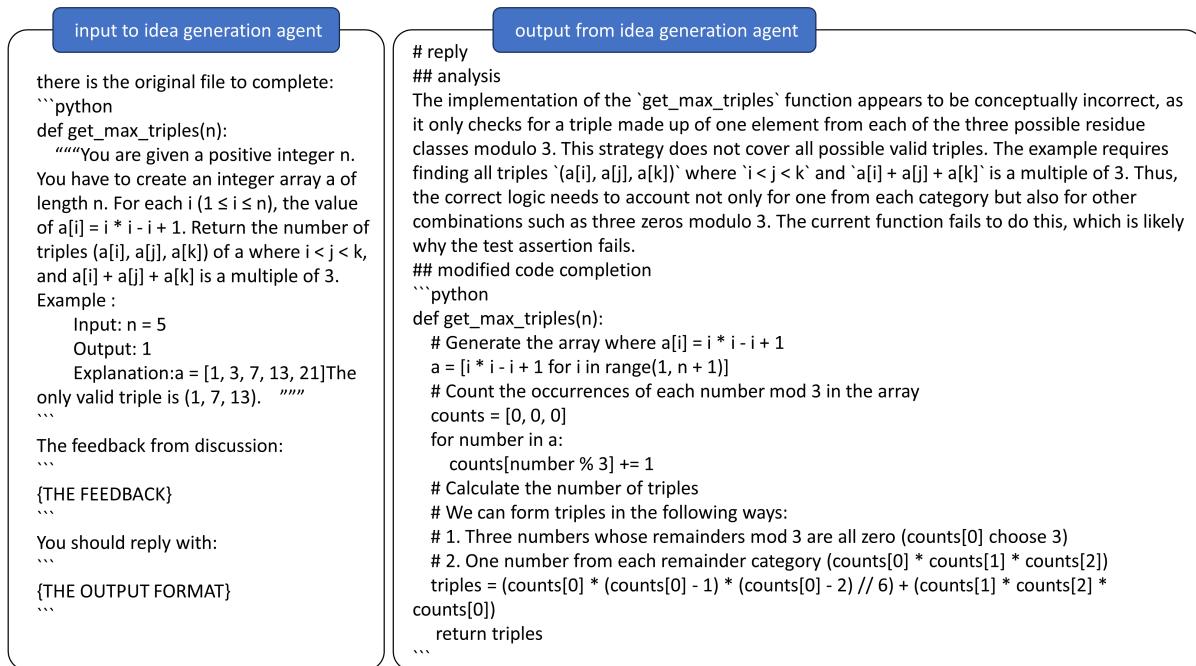


Figure 8: The original task is "You are given a positive integer n . You have to create an integer array a of length n . For each i ($1 \leq i \leq n$), the value of $a[i] = i \times i - i + 1$. Return the number of triples $(a[i], a[j], a[k])$ of a where $i < j < k$, and $a[i] + a[j] + a[k]$ is a multiple of 3. Example: Input: $n = 5$, Output: 1, Explanation: $a = [1, 3, 7, 13, 21]$, The only valid triple is $(1, 7, 13)$."

At first, one of the idea agents gives a draft saying that for all numbers in a_i , there is only 3 possible reminders when divided by 3 : 0,1,and 2. so it suggest first counts how many numbers reminds 0,1,and 2. than there should be 2 senario that the 3 selected 3 numbers from the list satisfy the requirement : all the reminder is 0 or reminder is 0,1,2 .these two can be calculated by the 3 counts .

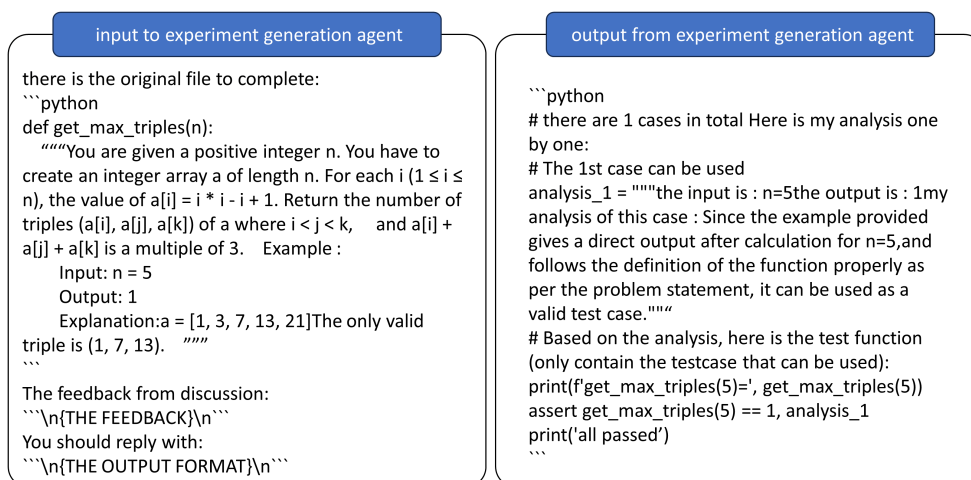


Figure 9: then , one of the experiments agents , design its experiment about the task , first print the result if $n=5$, then if the code draft showed that the result is not 5 , raise an error message contains the analysis about why when $n =5$ it should be the result .

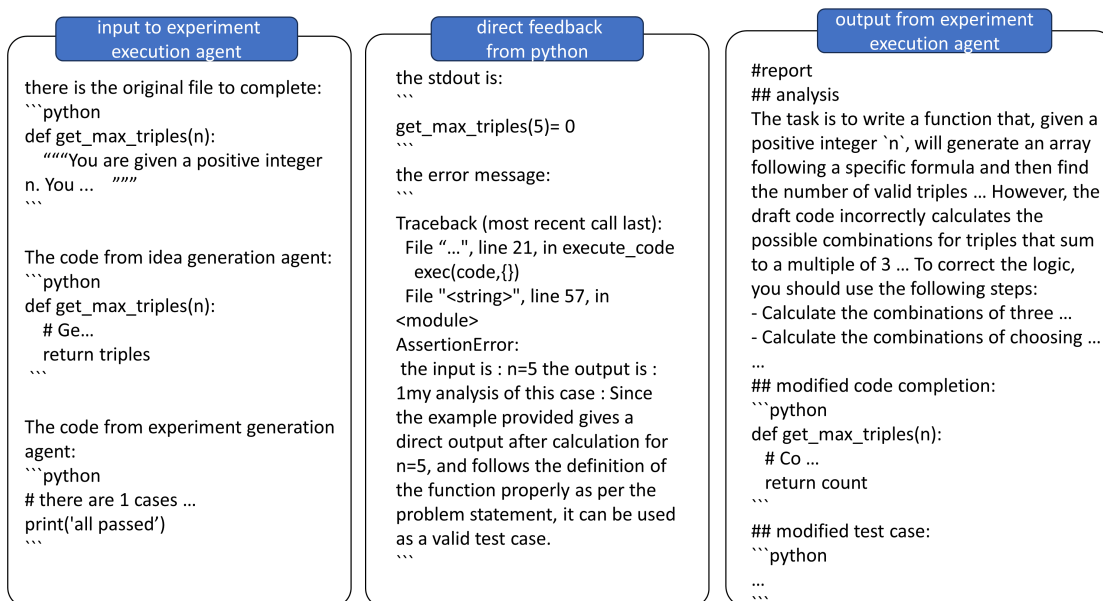


Figure 10: then the experiment executed in a python interpreter , and give stdout , error message , the high light part of the error message is added by the experiment agent . this part is augmented feedback of current code draft .

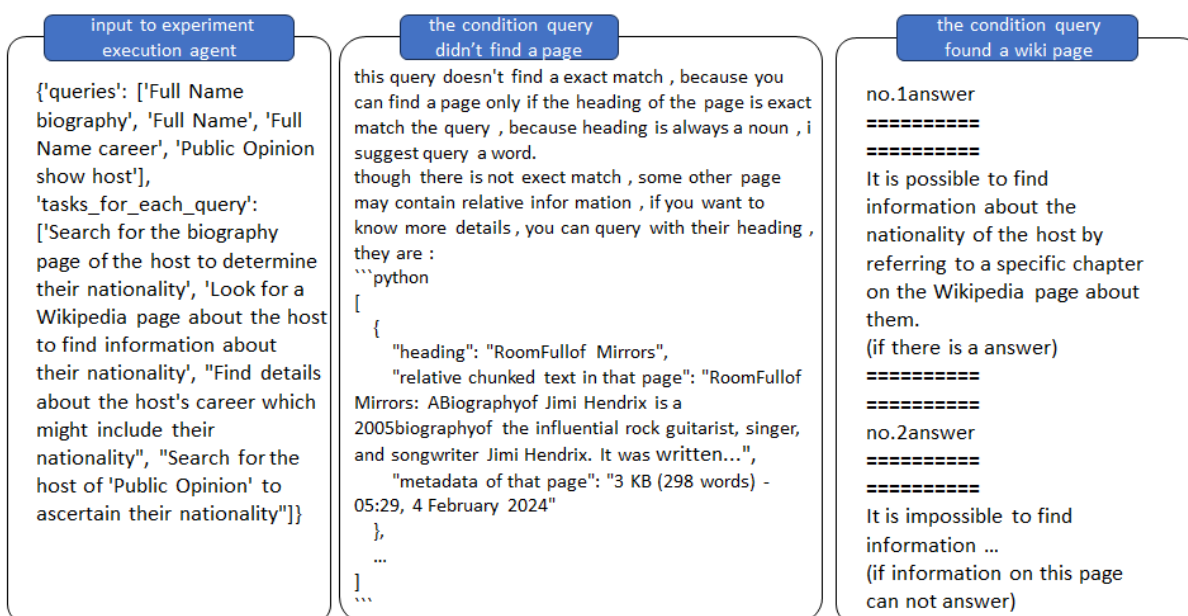


Figure 11: Input output example for experiment execution agents in HotpotQA

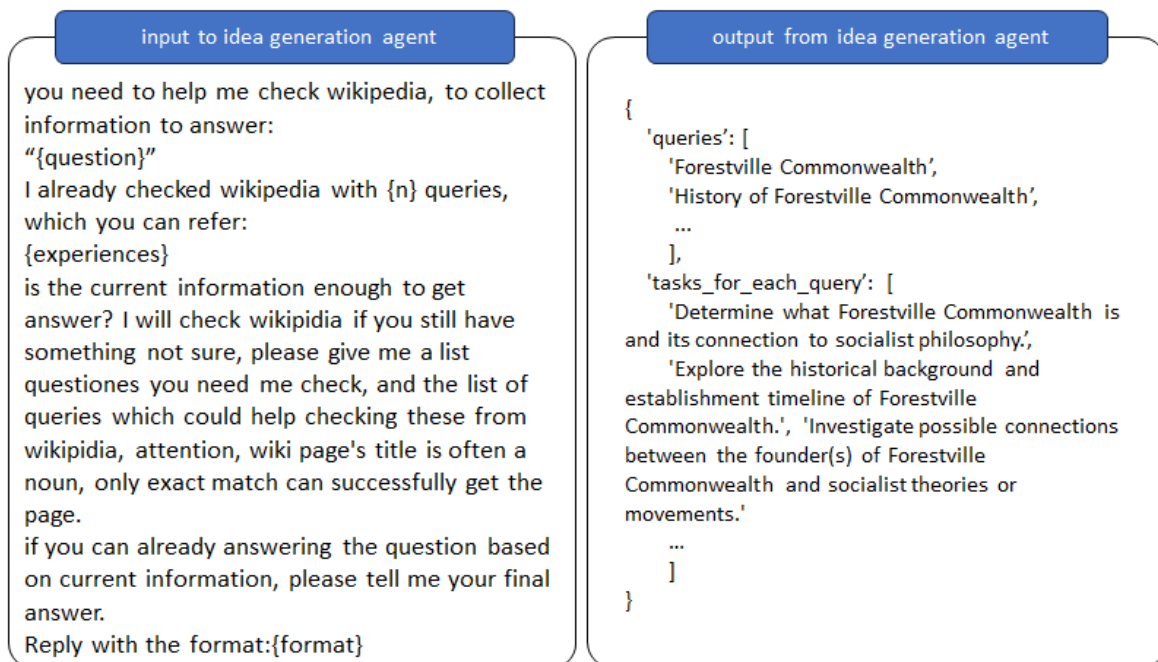


Figure 12: Input output example for idea generation agents in HotpotQA