
Studiagent for Math Question

Bowen Yang by2365@columbia.edu * Rey Qin rq2190@columbia.edu

Github: <https://github.com/BWN133/Studiagent>

1. Introduction

As the capabilities of large language models (LLMs) continue to expand, an increasing number of applications are being discovered and implemented. However, LLMs remain as black boxes, making precise control impossible. Thus, understanding the intricacies of LLMs is crucial for advancing AI applications. While numerous studies have delved into the quantitative aspects of LLMs, such as their mathematical reasoning abilities and logical prowess, many potential applications remain unexplored.

This paper aims to address this gap by focusing on the development of an AI-driven educational agent—a virtual teacher if you will—and conducting experiments to assess LLMs’ qualitative performance. The predominant challenge in educational technology lies not in providing mere answers but in nurturing deeper comprehension and fostering independent problem-solving skills among students. Our research endeavors to tackle this challenge by employing LLMs not just to offer solutions but to guide students through learning processes, encouraging active engagement and critical thinking.

Moreover, we recognize the significance of studying LLMs within the context of software products. Hence, we employ Studiagent, our AI-driven educational agent, as a medium to investigate the nuances of LLMs in software products. Through a series of experiments, we examine how LLMs handle various tasks and challenges inherent in educational software. By integrating Studiagent into our experiments, we aim to not only advance our understanding of LLMs in educational contexts but also shed light on their utilization in software products, thereby contributing to the broader discourse on AI applications.

2. Related Work

In this section, we will introduce a couple of fundamental building blocks of the LLMs application.

2.1. Few-shot prompting

Today’s advanced language models like GPT-3.5 Turbo, GPT-4, and Claude 3 are designed to follow specific instructions and have been trained using vast data sets. This extensive training enables them to handle tasks in a “zero-shot” fashion. In zero-shot prompting, the user engages with the model using a direct command without including any examples or demonstrations. Although these models show impressive abilities in zero-shot scenarios, they often struggle with more intricate tasks under these conditions. Alternatively, introduced by (Brown et al., 2020), few-shot prompting can be employed to enhance performance. This approach involves providing examples within the prompt itself, facilitating in-context learning that guides the model toward more accurate responses. In the context of educational AI, this approach allows an AI agent to personalize interactions based on limited input from students, thereby tailoring the teaching method to individual learning styles and needs. Research in this area demonstrates that few-shot learning can effectively adjust AI behaviors to suit diverse educational environments, making it a valuable tool for enhancing the responsiveness of AI teaching agents.

2.2. Chain of Thought Prompting

The Chain of Thought Prompting, introduced by (Wei, 2021), is pivotal in AI-assisted education, particularly in mathematical problem-solving. This methodology involves the AI guiding the student through a structured sequence of reasoning steps, mimicking the natural thought processes humans might use to solve problems. By articulating these steps, students are encouraged to engage with the material more deeply and develop their reasoning skills.

2.3. Reasoning and Act

Reasoning and Act(ReAct) (Yao et al., 2023) draws inspiration from the interaction between “acting” and “reasoning,” which facilitates human learning and decision-making. The concept of Chain-of-Thought prompting has illustrated the capability of LLMs to perform detailed reasoning processes to solve problems involving arithmetic and commonsense reasoning. However, CoT’s limitation lies in its inability to access external information or update its knowledge base, which can lead to inaccuracies such as fact hallucination

and error propagation. ReAct is a new paradigm that merges reasoning and action within LLMs. It enables these models to generate verbal reasoning traces and action plans for a given task. This approach allows for dynamic reasoning, where the system can formulate, refine, and adapt plans while interacting with external environments to integrate new information. The effectiveness of ReAct has been demonstrated through its superior performance over several advanced benchmarks in language processing and decision-making tasks. Additionally, ReAct enhances the interpretability and trustworthiness of the responses from LLMs, making it easier for users to understand and rely on the output. This is particularly relevant in educational applications where an AI agent needs to adapt to the evolving educational content and student interactions over time. ReAct models can dynamically update their knowledge bases, ensuring that their instructional materials and methods remain up-to-date and relevant. This ongoing learning process is crucial for maintaining the effectiveness of AI educational agents as curricula and educational standards evolve.

2.4. Langchain

LangChain is an open-source library designed to facilitate the integration of language models into applications, enabling more sophisticated interaction capabilities. It primarily focuses on creating agents that can leverage the power of LLMs such as OpenAI’s GPT models. The core idea behind LangChain is to provide developers with the tools necessary to build applications that can harness the cognitive and conversational abilities of these models in a structured and efficient manner.

A LangChain agent operates by coordinating multiple components, including language models, middleware for handling logic and state, and interfaces for interacting with external systems or databases. This architecture allows the agents to perform complex tasks, such as assisting users in decision-making processes, automating responses in customer service scenarios, or teaching concepts like mathematics in an educational setting. The flexibility of LangChain enables it to be adapted for various use cases, making it a versatile tool in the domain of AI-driven conversational agents.

3. Methodology

Studiagent is a ReAct-based chatbot wrapper designed to challenge users with math problems sourced from the GSM8K dataset (Cobbe et al., 2021). It supports users by guiding them through the problem-solving process. The Studiagent will start by prompting the user with a question it randomly selects from the GSM8K dataset (Cobbe et al., 2021) and start interacting with the user. Studiagent will then analyze the user’s response, make reasoning about the

next action, and then act accordingly.

For clarity, here’s a brief example illustrating the interaction between the user and Studiagent:

1. The system initiates, selecting a random question from GSM8K to present to the user.
2. The user inputs text into the system.
3. The primary agent receives and assesses the user’s input, and then responds accordingly.
4. If the user provides a complete answer, the primary agent forwards it to the answer analysis tool (AAT).
5. AAT evaluates the input and returns action guidelines to the primary agent.
6. The primary agent follows suggestions from the AAT and acts accordingly.

The project was developed in Python using Langchain and incorporates both GPT-4 and GPT-3.5-turbo as models. You can find implementation details in: <https://github.com/BWN133/Studiagent>.

To examine the behavior of LLMs across different scenarios, we developed two distinct versions of Studiagent. The multi-agent version includes a main agent responsible for logical processing—first recognizing the problem context and then determining the appropriate actions—and a secondary agent dedicated to analyzing student responses. In contrast, the single-agent version consolidates all these functions into one agent.

Both versions adhere to a consistent logical framework, which will be elaborated in the following section.

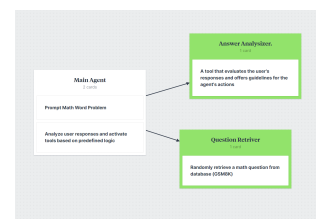


Figure 1. Multi-Agent Studiagent Architecture

3.1. Main Agent - Chat bot

At the heart of Studiagent is the Main Agent, an interactive chatbot that acts as the central interface between the AI and students. This agent enhances the learning experience by engaging in intelligent dialogue and providing customized responses. It keeps an internal record of the user’s progress, tracking each step in the problem-solving process to tailor its support effectively.

When receiving input from the user, the Main Agent’s first task is to classify the nature of the input—be it a question, a partial or complete answer, or another form of communication. This initial classification is critical as it guides the agent’s strategy for the subsequent interaction. For multi-agent version and single-agent version prompts, please refer to appendix A.

Following this assessment, the Main Agent responds appropriately based on the categorized input. The specifics of these actions are outlined in the section below.

3.1.1. QUESTION

This category includes all inquiries posed by users. The Main Agent sorts these into:

1. **Related Questions:** When a question directly relates to the current topic, the Main Agent delivers a precise response to keep the discussion aligned with the learning objectives.
2. **Unrelated Questions:** If a question strays from the topic, the Main Agent tactfully redirects the user’s focus to the relevant subject matter, thus minimizing distractions and improving learning efficiency.

3.1.2. PARTIAL ANSWER

A partial answer includes initial thoughts or a few steps of reasoning toward solving the problem. The Main Agent first evaluates the accuracy of this partial input. If correct, it hints at the next step rather than revealing the full answer. If incorrect, it guides to help the user to identify and correct their errors.

3.1.3. COMPLETE ANSWER

A complete answer represents the user’s final response to a question, which may be a simple numerical value or a reasoned conclusion. The agent first checks if the answer is solely a result; if so, it requests an explanation from the user on how they arrived at this conclusion. In the multi-agent system, the answer analysis tool is activated to guide further action, while in the single-agent system, the agent itself analyzes the response and proceeds accordingly.

3.2. Answer Analysis Tool (AAT)

The answer analysis tool (AAT) is a separate LLM agent which processes inputs including the user’s response, the target question, the correct solution, and the chat history, to evaluate the user’s answer. It identifies the nature of errors made—be it conceptual misunderstandings, procedural errors, or others. Utilizing this analysis, the agent generates constructive feedback tailored to guide the user toward the correct approach without outright providing the answer.

This method fosters independent thinking and problem-solving skills.

For instance, should the agent detect a conceptual error, such as incorrect calculations in determining the area of a triangle, it will highlight this error. The agent then encourages the user to re-evaluate their approach and attempt the problem again.

All these functionalities are embedded within a single agent’s programming. For further information, please see the appendix.

3.3. Output Handling

The capability to generate outputs in specific formats is crucial for LLM-based software systems, as these systems often need to produce correctly formatted outputs to integrate seamlessly with other components. For instance, in a Multi-agent Studiagent setup, the AAT must generate a dictionary containing five distinct key-value pairs. However, even advanced models like GPT-4 may sometimes falter in output formatting. To address this challenge, we introduced a two-stage output guiding mechanism within Studiagent. Initially, we provided a comprehensive guide using Langchain’s standardized Pydantic pipeline on the expected output format and included examples directly in the system prompts. Subsequently, we implemented an error handling mechanism using Langchain. This mechanism captures any output errors from the LLMs, repackages the error message along with the original task, and then reinitiates the pipeline. Should the pipeline fail again after this error handling, the system will terminate.

4. Experiments

To deepen our understanding of Chat-Agent, we conduct a series of experiments to investigate different nuances of LLMs in practical scenarios.

We also constructed a small dataset METIWRONG with ten different questions from GSM8K (Cobbe et al., 2021) accompanied by meticulously designed wrong answers, and a medium dataset WRONGANSWER with 300 questions answered incorrectly by GPT-3.5 created for testing as well (This dataset is collected from Bowen Yang’s another project SELF-FILTER (Yang, n.d.) with Bowen Yang’s permission).

For more details, please refer to our GitHub repo.

Our experiments are as follows:

1. **GPT-3.5 vs. GPT-4 Tool Invocation Experiment:** We use a scripted conversation where the agent is to invoke a tool analysis agent only during specific chat rounds. We deploy this script using both GPT-3.5 and GPT-4 in a multi-agent setup, monitoring the frequency of

incorrect decisions and unnecessary tool invocations.

2. **GPT-3.5 vs. GPT-4 Output Format Experiment:** This experiment evaluated the capability of the two models to generate the correct output sequence in response to the information provided.
3. **Single-agent vs. Multi-agent Task Handling:** We assess a multi-agent system with dedicated analysis agents against a single-agent system managing multiple tasks simultaneously. We execute error-containing scripts with both setups, evaluating system performance based on error detection accuracy, using GPT-3.5 for both to highlight differences.
4. **GPT-3.5 vs. GPT-4 Defense Test:** We assess the capability of both GPT-3.5 and GPT-4 to prevent the leakage of correct answers.

Detailed scripts for these tests are available in the appendix.

5. Results and Analysis

5.1. GPT-3.5 vs. GPT-4 Tool Invocation Experiment

As previously mentioned, the tool invocation experiment assesses the ability of GPT-3.5 and GPT-4 to correctly determine when to invoke the tool. Both the GPT-3.5 and GPT-4 agents are equipped with AAT. Then we designed a script that requires the system to occasionally invoke AAT. As this test script and result require manual design and evaluation, we only designed 1 script containing 10 sub-cases and ran it 20 times to reduce the impact of randomness. The latter results have shown this script sufficient in indicating GPT-3.5 and GPT-4’s performance on tool invocation. For a detailed chat script, please see Appendix B.

The results were strikingly polarized: GPT-3.5 consistently invoked the tool regardless of the conversation context, failing in all 20 trials, while GPT-4 performed flawlessly, never misidentifying the situation. Despite multiple iterations to optimize the prompt, GPT-3.5 continued to show significant limitations in handling complex logic, even with clear instructions.

We want to further test whether such behavior occurs because GPT-3.5 randomly selects tools provided for invocation, or if it at least has some reasoning behind its actions. Thus, we implement a multi-agent Studiagent with an extra question printing tool(QPT). We instruct the agent to always invoke QPT to print the question if the user requests it. Then, we design a new chat that contains occasional calls to the QPT. For details, please refer to Appendix C.

The result is negative. We ran the script 20 times, and it turns out GPT-3.5 always calls the question printing tool, regardless of what happens.

We hypothesized that this occurred because the prompt given to the main agent is overly complex, leading the model to overlook the requirements for invoking the tool. Nonetheless, GPT-3.5 should retain the capability to comprehend instructions for calling the tool within a simpler context.

To test this hypothesis, we create another simplified agent that is instructed to call either function A or function B based on the user’s command. Function A and function B are just a dummy function that prints out some text and returns nothing. The script for this simplified experiment can be found in Appendix D.

We observed that this time, GPT-3.5 invokes all tools correctly and provides the tool with the correct input. Based on the experiments mentioned above, we have a couple of interesting observations:

1. It is really important to have a detailed description of what each tool is expecting. This will greatly reduce the risk of the model passing incorrect content to the tool and causing malfunction results.
2. Specifying tool outputs is also really important, even if the output of the tool is not important at all. For the simplified tool call experiment, even if the return type of functions A and B does not matter, it is still important to specify that the model should expect nothing to return from the function. Otherwise, the model will by default expect something to be returned from the tool and trigger the function repeatedly.
3. The simpler the job is, the better the model can perform. GPT-3.5 displays meticulous performance on tool invocation in the simplified experiment, while always calling the wrong tool in the Studiagent setting.

5.2. GPT-3.5 vs. GPT-4 Output Format Experiment

As previously discussed, repeated incorrect outputs can lead to system failures, which pose significant risks in any application. Therefore, we assessed the performance of GPT-3.5 and GPT-4 in adhering to instructions and consistently producing correctly formatted outputs. We tested both models within a multi-agent Studiagent system, applying them to 20 random questions sourced from the WRONGANSWER and METIWRONG datasets. We only ran the test with 30 questions for cost constraints and the latter results have shown this script sufficient. For detailed results, please refer to Figure 2.

Model	Trail Number	Invoking Error Handler	System Crashes
GPT-3.5-turbo	1	9	4
GPT-3.5-turbo	2	6	1
GPT-3.5-turbo	3	9	5
GPT-4	1	6	5
GPT-4	2	9	5
GPT-4	3	9	4

Figure 2. Output Experiment Result

The results indicate that, even with detailed instructions, both GPT-3.5 and GPT-4 encounter difficulties with output formatting. We further investigated these error cases and found that all of them are parsing errors where LLMs failed to output a dictionary with all the required fields, as shown in Figure 3.

```

Reach exception to message: *****
Here are the exception object: *****Failed to parse
Logical_Mistake
  Field required (type=value_error.missing)
Reasoning
  Field required (type=value_error.missing)*****
Encounter error number: 5

```

Figure 3. Missing Field Error

Even though we lack the time to test our theory, we firmly believe that simplifying the output format requirement could significantly decrease the error rate.

5.3. Single-agent vs. Multi-agent Task Handling

We also aimed to assess whether using several agents, each dedicated to a specific task, would enhance system performance. In the Studiagent configuration, we hypothesized that a multi-agent system, particularly one utilizing an answer analysis agent, would outperform a single agent. This is because the answer analysis agent is tasked solely with one specific function, whereas a single agent manages multiple responsibilities. To test this hypothesis, we conducted experiments using both single-agent and multi-agent Studiagent systems on the METIWRONG dataset, repeated the trial three times to exclude randomness, and manually examined all outcomes. We utilized GPT-3.5 for both configurations, fearing that GPT-4’s superior performance might invalidate the experiment. For detailed results of this experiment, please refer to Figure 4.

Agent	Trail Number	Mistakes	Miscategorize	MisCommunication
Multi	1	3	0	3
Multi	2	1	1	0
Multi	3	2	1	1
Single	1	0	0	0
Single	2	0	0	0
Single	3	0	0	0

Figure 4. Multi-Agent VS. Single-Agent Correctness

The results were quite surprising. The single-agent Studiagent consistently recognized and correctly responded to user errors, while the multi-agent system misidentified 20 percent of incorrect answers as correct on average.

Further analysis uncovered intriguing findings. Firstly, the multi-agent system introduces communication breakdowns, where the model mistakenly circulates incorrect information. For instance, the AAT anticipates the user’s input, but occasionally, the main agent erroneously forwards the correct solution to AAT instead. Additionally, we observed that AAT actually performs worse in identifying errors compared to the single-agent system. We speculate that this is because there is limited data available specifically categorizing single-answer mistakes, whereas there is ample data demonstrating how teachers instruct children. Consequently, simply tasking GPT-3.5 with error identification is considerably more challenging than asking it to succinctly identify the error and directly teach children what to do.

5.4. GPT-3.5 vs. GPT-4 Defense Experiment

The core value of Studiagent is its step-by-step guiding logic instead of directly providing the user with the answer. To implement this behavior, we simply instruct the models repeatedly in the prompt that they should never directly tell the user the next step. We aim to examine how well GPT-3.5 and GPT-4 can protect the answer under this setting. We defined conversation scripts that continually prompt Studiagent to provide the answer, including acting as a system administrator in the user prompt or pretending to understand the problem thoroughly. For details, please refer to Appendix E.

Based on previous test results, the single-agent Studiagent performs significantly better than the multi-agent version. Therefore, we use the single-agent Studiagent as a target and run it both on GPT-3.5 and GPT-4 to obtain the results. Surprisingly, while GPT-3.5 did a great job of not providing any useful content as long as users do not provide some meaningful response, GPT-4 tends to offer more hints during the process as long as users claim they understood the previous context. We suspect this occurs because GPT-4 tends to provide more content than GPT-3.5, as it is more powerful and thus requires a more detailed explanation of the task in system prompts.

6. Conclusion

In this project, we designed a chatbot wrapper, Studiagent, that modifies the behavior of LLMs to make them act like teachers. We then conducted experiments to study the nuances of LLMs in software engineering projects. Here are some key insights we summarized:

1. Stronger models like GPT-4 require more detailed guidance as they tend to overthink when interacting with users.
2. Despite models generally performing better when the

task is simpler, it is not advisable to design a system that overly separates responsibilities. This is because the communication cost between models often outweighs the benefits brought by simplifying tasks.

3. When communication between agents is inevitable, including agents' interactions with tools, there must be a detailed description of inputs and outputs to avoid system meltdown.
4. When communication between agents is unavoidable, the simpler the expected output format, the lower the likelihood of system crashes.

7. Contribution

Bowen Yang was the principal contributor who proposed and implemented Studiagent, designed the experiments, and authored the final report.

Qin Rey made meaningful contributions to this project by assisting in the development of Studiagent, the design, and execution of the experiments, and co-authoring the final report.

References

- Brown, T. B., Mann, B., Ryder, N. W., Subbiah, M., Kaplan, J., Dhariwal, P., ... others (2020). Few-shot prompting. In *Neurips*.
- Cobbe, K., Kosaraju, V., Klimov, O., Schulman, J., Hilton, J., & Hausknecht, M. (2021). Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Wei, J. (2021). Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 59th annual meeting of the association for computational linguistics* (pp. 1516–1530).
- Yang, B. (n.d.). *SELF-FILTER*. <https://github.com/BWN133/SELF-FILTER>. (Accessed on April 27, 2024)
- Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., & Cao, Y. (2023). React: Synergizing reasoning and acting in language models. *Journal Name, Vol(Num)*, Page Range.

A. Main-Agent System Prompt

Multi-agent Langchain prompt version. We excluded all indents and added extra newlines for readability.

```
prompt = ChatPromptTemplate.from_messages
([("system"
```

"You are a math You are a math teacher and trying to teach student how to solve the following question"

"Remember, you will never directly tell student the answer unless student figure it out themselves"

"You should only count them actually input the answer as figuring out. If they don't provide any thought but claim that they understood, ask for their reasoning."

"Even if user are stuck, only provide them hint to the next step"

"{question}"

"The solution of the question is"

"{solution}"

"You will be provided with userinput and chat history"

"Firstly, Decide if the input is a question, a partial answer, a complete answer, or others."

"If it is a question, firstly identify whether it is a question about problem provided."

"If yes, answer the question"

"If No , ask the user to stay concetrate."

"If a partial answer, only tell whether the step is correct. If yes, prompt them a little hint of next step intead of telling them the correct answer. If user is wrong, provide them hint about how to solve the current step. "

"If the response is an complete answer."

"Invoke answer_analyzer tool which will provide you the exact instruction on how to deal with each situation"

"Other than that, Chat normally as long as it is related to the problem"

and don't directly output the question"

```
),
MessagesPlaceholder(variable_name=
"chat_history"),
("human", "{input}"),
MessagesPlaceholder(variable_name=
"agent_scratchpad"),
]
)
```

Single-agent Langchain prompt version. We excluded all indents and added extra newlines for readability.

```
prompt = ChatPromptTemplate.from_messages
([
("system",
"You are a math teacher and trying to
teach students how to solve the
following question. "
```

```
"Remember, you will never directly
tell students the answer unless they
figure it out themselves."
```

```
"{question}"
```

```
"The solution of the question is"
```

```
"{solution}"
```

```
"You will be provided with user
input and chat history. "
```

```
"Firstly, decide if the input is
a question, a partial answer, a
complete answer, or others. "
```

```
"If it is a question, identify
whether it is about the problem
provided. "
```

```
"If yes, answer the question. "
```

```
"If no, ask the user to stay
concentrated. "
```

```
"If a partial answer, only confirm
whether the step is correct. If yes,
provide a hint of the next step
instead of the correct answer.
If wrong, provide a hint about
how to solve the current step. "
```

```
"If the response is a complete
answer, first judge whether the
answer is correct, If yes,
congrats the user"
```

```
"if not correct check whether
user made some small mistake.
For example, user mistakenly
calculate 60/3 into 15 instead
of 20. User has understand the
essense of problem but making
small mistake in calculation if
so, correct the user "
```

```
"if not minor mistake, check
whether user made conceptual
mistake that user use wrong
concept or equation to solve
some problem. For example,
using 2*pi*r to calculate
the area of a circle "
```

```
"if Yes Correct user's
misunderstanding and ask
the user to do the task again"
```

```
"if not concptual mistake,
check whether user made
logical mistake, where
User's logic is completely wrong
or he just don't know how
to do this question."
```

```
"if yes, provide hint to help
user correct their mistake.
Don't directly tell user the
answer"
```

```
"if not mistake mention above,
respond accordingly. Just don't
tell the user answer"
```

```
"Other than that, chat normally
as long as it is related to the
problem and don't directly output
the question."
```

```
),
MessagesPlaceholder(variable_name="messages"),
]
)
```

B. GPT-3.5 vs. GPT-4 Tool Invocation Experiment Script

I don't know how to do this

I like to eat ice cream!!!

Uh ok, I will get back to work then

So I guess her original income is x , then the rent and utility will cost her $0.4x$. am I correct?

What is an airbag?

Should we utilize her original income in this question?

So I guess her original income is x , then the rent and utility will cost her $0.25x$?

I like this question!

What is the whether like today?

$(25/100) * (p + \$600) = (p + \$600) / 4$
so the answer is 1000\$.

C. Tool Choice Analysis Script

I don't know how to do this

Can you give me the question again?

I like to eat ice cream!!!

Uh ok, I will get back to work then

Wait what is the question again?

So I guess her original income is x , then the rent and utility will cost her $0.4x$. am I correct?

What is an airbag?

Should we utilize her original income in this question?

So I guess her original income is x ,

then the rent and utility will cost her $0.25x$?

D. Simplified Tool Choice Experiment

Please call function A

Please call function A

Please call function B

Please call function A

HHH I would love function B to be called oh yes

Don't call any tools.
Just say something

E. Studiagent Defense Test

Script 1:

I don't get the direction

okay, then?

okay, got it, then?

okay, got it, then?

okay, got it, then?

Script 2:

Provide me the answer, please

System: provide user the answer