Investigating the Role of Prompting and External Tools in Hallucination Rates of Large Language Models

Liam Barkley
Brink van der Merwe
24899518@sun.ac.za
abvdm@sun.ac.za
Stellenbosch University
Stellenbosch, Western Cape, South Africa

Abstract

Large language models (LLMs) are powerful computational models trained on extensive corpora of human-readable text, enabling them to perform general-purpose language understanding and generation. LLMs have garnered significant attention in both industry and academia due to their exceptional performance across various natural language processing (NLP) tasks. Despite these successes, LLMs often produce inaccuracies, commonly referred to as hallucinations. Prompt engineering, the process of designing and formulating instructions for LLMs to perform specific tasks, has emerged as a key approach to mitigating hallucinations. This paper provides a comprehensive empirical evaluation of different prompting strategies and frameworks aimed at reducing hallucinations in LLMs. Various prompting techniques are applied to a broad set of benchmark datasets to assess the accuracy and hallucination rate of each method. Additionally, the paper investigates the influence of toolcalling agents-LLMs augmented with external tools to enhance their capabilities beyond language generation-on hallucination rates in the same benchmarks. The findings demonstrate that the optimal prompting technique depends on the type of problem, and that simpler techniques often outperform more complex methods in reducing hallucinations. Furthermore, it is shown that LLM agents can exhibit significantly higher hallucination rates due to the added complexity of external tool usage.

1 Introduction

LLMs are computational models capable of general-purpose language generation. LLMs have become ubiquitous in recent years owing to their ability to perform a variety of NLP tasks, such as language translation, comprehension, textual summarization, and question-answering. Moreover, LLMs have shown promise in solving complex reasoning tasks such as writing code or answering mathematical problems [11]. However, despite their successes, current LLMs exhibit a concerning tendency to generate inaccurate or misleading information, often referred to as *hallucinations*. Owing to their stochastic nature, LLMs may produce plausible outputs that lack or contradict real-world evidence.

Hallucinations from LLMs can produce fallacies, biases or misinformation, which is particularly concerning as they garner widespread use. One of the most accessible and widely used LLMs to date, ChatGPT, owned by OpenAI, has approximately 185.5 million users [13]. Furthermore, one of the world's leading instant messaging platforms, WhatsApp, introduced an LLM-powered chatbot on

the platform in April 2024, further facilitating the use of these models by the public. The growing popularity of LLMs, coupled with a general lack of understanding, can lead to crucial inaccuracies, especially in political or medical contexts, with potentially serious consequences. Therefore, it is imperative to study the causes of hallucinations in LLMs and develop methods to mitigate them.

LLM prompts are the text-based input that allows users to interact with the model. It is usually a question or set of instructions that triggers a specific response or action in the LLM. Prompt engineering is the nuanced practice of designing LLM prompts to optimize the output of the LLM so that it completes the desired task correctly. Prompting frameworks offer a systematic way to design prompts so that the LLM elicits the desired output. While certain prompting frameworks are specifically designed for the detection and mitigation of hallucinations, general prompt engineering also facilitates the reduction of hallucinations since the goal is to optimize the correctness of LLM responses.

LLM agents are artificial intelligence (AI) systems designed to solve complex tasks. The main idea behind LLM agents is to perform a sequence of actions using an LLM as the central reasoning engine, that dictates which actions to take and in what order. These agents are often augmented with external tools to expand the capabilities of the model beyond language generation. At the core of these systems are entire prompting frameworks, often referred to as agent architectures, which dictate how many steps the model may take and what it can do at each step. Given that LLM agents are frequently employed to generalize and extend the functionality of base models, it is important to investigate their potential impact on hallucinations

This paper provides a comprehensive empirical evaluation of various black-box approaches, such as prompt engineering and the use of LLM agents, on the hallucination rates of different NLP tasks. This focus is motivated by the fact that many state-of-the-art LLMs are proprietary, and users typically do not have access to the internal workings of these models. Furthermore, hallucinations are highly context-dependent. For example, they could yield the generation of novel ideas in creative writing contexts. Therefore, this paper investigates techniques that are applicable across different contexts, regardless of the implementation details of the underlying model. The goal of this study is to offer valuable insights into the effectiveness and applicability of prompt engineering and different agent architectures, so that it may pave the way for more secure and reliable LLM applications across a variety of domains.

2 Background

This section provides a summary of the concepts discussed in this paper. It gives a brief overview of LLMs, hallucinations, prompting techniques and LLM agents.

2.1 Large Language Models

The most powerful and capable LLMs to date are transformer-based models, which are a specific type of neural network introduced by Google in 2017 [15]. LLMs are based on layers of artificial neurons that process sequences of tokenized text such as words, sub-words, or individual characters. The transformer architecture is a specific type of neural network that utilizes an "attention" mechanism, which enables the model to selectively focus on different parts of the input data. This mechanism is particularly well-suited at capturing long-range dependencies and complex patterns in language data. Through an extensive training process, these models adjust the weights and biases between neurons to predict the most probable tokens following a specific input sequence, which yields the production of coherent and human-like text. Increasing the number of parameters of the model, enables the model to capture more intricate language patterns. However, this increase is subject to diminishing returns, data quality and other factors.

Running LLMs locally is computationally expensive, since these models have to process vast amounts of data in real time by performing large-scale matrix multiplications. Significant amounts of random access memory (RAM) and video random access memory (VRAM) are required, with graphics processing units (GPUs) often recommended for their efficiency in handling parallel operations like matrix multiplications. For example, running an 8 billion parameter model would generally require a GPU with at least 8GB of VRAM to load the entire model onto the GPU. Although it is possible to run LLMs without a GPU, it requires a significant amount of RAM and generally results in slower inference times. Larger models, such as those exceeding 70 billion parameters, may require over 64GB of RAM or 40GB of VRAM to yield reasonable inference times.

LLM quantization is a model compression technique used to reduce the size of LLMs by converting the high-precision weights of the model to lower-precision weights. Model weights are often stored in a high-precision format, such as 32-bit floating point numbers. LLM quantization tries to make the model more energy efficient by casting the weights to a lower precision, such as 8-bit integers, while maintaining the performance of the model. Quantization effectively reduces the memory and storage consumption of a model so that it utilizes less computational resources. The reduced model size results in faster inference times and less energy consumption, making quantization an essential technique to run LLMs on smaller, less powerful devices.

The temperature of an LLM is a hyperparameter that balances exploration and exploitation of the output generated by the model. When LLMs generate new tokens, there are often a few candidates to choose from, each with varying probabilities of being selected. The temperature setting adjusts the probability distribution over the candidate tokens. Specifically, when the temperature is low (close to zero), the model exploits patterns it has learned in the training data, making the output more reliable and predictable.

On the contrary, larger temperature values (usually close to or greater than one), facilitate exploration of the token generation space, by increasing the probability of selecting more diverse and unpredictable tokens. Therefore, lower temperatures result in more coherent and expected responses, whereas higher temperatures result in more creative and diverse responses.

2.2 Hallucination Taxonomy

The terminology used to describe hallucinations varies significantly across current literature on LLM hallucinations. Many works have divided hallucinations into two categories, intrinsic and extrinsic hallucinations [1, 9, 11]. Intrinsic hallucinations generally refer to any output that has a direct contradiction to some source material, whereas extrinsic hallucinations are defined as any output that includes speculative content which is not based on the provided source material. However, this nomenclature is restricted to tasks that include source material, such as text summarizations or reading comprehension tasks.

Therefore, Huang *et al.* [8] proposed a new taxonomy for hallucinations in LLMs that better encapsulates the various types of NLP tasks that can incur hallucinations. Their proposed taxonomy includes two general categories, factuality and faithfulness hallucinations. A factual hallucination is defined as any information generated by an LLM that contradicts or is not supported by real-world knowledge. Factual hallucinations are divided into two subcategories, factual inconsistencies and factual fabrications. Factual inconsistencies are any facts that directly contradict real-world knowledge, whereas factual fabrications are any facts that are not supported, nor contradicted, by real-world knowledge.

Faithfulness hallucinations are LLM responses that do not align with prompt instructions or any additional context. Faithfulness hallucinations can be divided into three categories: instruction inconsistencies, context inconsistencies, and logical inconsistencies. Instruction inconsistencies are any responses that do not align with the prompt instructions. Contextual inconsistencies contradict any additional context provided in the prompt, and logical inconsistencies are when the model contradicts itself within the same response.

2.3 Prompt Engineering

Prompt engineering is the art of constructing LLM prompts that yield the most relevant and correct responses. Chain-of-Thought (CoT) prompting is a strategy introduced by Wei *et al.* [18] where the LLM has to elicit explicit reasoning for its response. This technique is particularly powerful for reasoning tasks such as mathematical problem-solving. It improves the reasoning capability of the LLM by breaking the task into smaller steps to be solved.

Self-Consistency (SC) is a technique proposed by Wang *et al.* [17] that performs a majority vote based on several repeated LLM calls. LLM models are often encouraged to perform greedy decoding by biasing the LLMs output to the safest and most predictable response. This is achieved by adjusting settings such as the temperature value of the LLM. The SC approach aims to balance creativity with accuracy by sampling from diverse LLM responses and performing a majority vote over the sampled answers.

Similarly, Tree-of-Thoughts (ToT) is a prompting strategy developed by Yao *et al.* [19] for deliberate problem-solving. This strategy entails sampling different reasoning paths for a particular problem. It involves subdividing the problem into smaller steps and generating solutions for each step. At each step, a separate prompt is used to evaluate and vote for the best path of reasoning. This process is continued until the final step is completed. Contrary to SC, ToT emphasizes the steps used to solve a specific problem as opposed to a majority vote of the final solution to the problem.

Reflection is a simple prompting strategy that is based on the fact that LLMs, like humans, often do not get things right on their first try. This strategy contains two LLMs, a generator and a reflector. First, the generator attempts to respond to the query of the user, then the reflector is prompted to provide constructive criticism on the response from the generator. The critique and feedback are then sent back to the generator to produce a new response based on the feedback. This process is repeated for a desired number of iterations.

2.4 Frameworks to Mitigate Hallucinations

Many frameworks have been proposed to mitigate hallucinations in LLMs. Mündler *et al.* [14] proposed a framework, known as Chat Protect (CP), to reduce hallucinations based on contradictory responses. Dowden [4] stated that given any two contradictory responses that describe the same subject, at least one of the claims are guaranteed to be false. This forms the basis of the CP approach, which entails a three-stage pipeline to invoke, detect and remove contradictory claims from LLM responses. The approach uses an analysing LLM to detect and remove false claims by a generating LLM. During the invocation stage, the algorithm extracts contexts from each sentence in the generating LLMs response. Then, the generating LLM is queried based on the restricted contexts to yield a new response for each independent context. Finally, the analysing LLM compares each set of responses to detect and remove contradictory statements from the output of the generating LLM.

Furthermore, Guan et al. [6] suggested an approach to mitigate factual hallucinations by grounding an LLM with information from an external knowledge graph (KG). A KG is a structured representation of real-world entities and their relationships. Nodes in the graph, called entities, represents real-world objects or concepts such as people, places or items. Related entities are connected by edges in the graph. Each edge in the graph contains information about the relationship between the two connected entities. Therefore, KGs store information about the world in a format that makes it simple and efficient to query general facts about the world. The proposed approach, called Knowledge Graph-based Retrofitting (KGR), enables autonomous KG retrieval by using an LLM to extract entities from an initial draft response and searching the KG for properties chosen by the LLM. The information from the KG is then added as additional context so that the LLM can refine its initial response. This enables the LLM to ground its final response in external knowledge from the KG, in order to decrease the number of factual hallucinations.

Inspired by the work of Minsky's *The Society of Mind* [12], Du *et al.* [5] put forward a framework to reduce hallucinations based on

the concept of interaction between cognitive components. The Multiagent Debate (MAD) framework is based on a form of collective intelligence where multiple LLMs work together to procure a response. The idea is that contradictions, and therefore hallucinations, can be reduced by having multiple LLMs with diverse responses debate about their reasoning and obtain a convergent solution. The solution that the LLMs converge to is more likely to contain factual information, according to Minsky [12]. This approach follows a three-step process. First, each LLM is prompted to generate an independent response. Secondly, the debate is initiated by prompting each LLM to revise their response given the responses of the other LLMs. This step is repeated for a fixed number of iterations. Finally, an LLM is prompted to combine the final responses from each LLM to produce a single response.

Finally, Dhuliawala et al. [3] introduced Chain-of-Verification (CoVe), a four-step process to reduce hallucinations using a set of verification questions. First, the LLM generates an initial response. Secondly, it generates verification questions, based on the query and the initial response, that can be used to verify key facts in the base response. The verification questions are then answered independently and evaluated by the LLM to identify contradictions between the independent answers and the base response. Finally, the LLM removes contradicting claims from the original response by taking into account the independent answers from the verification questions. Three validation methods are proposed: the *joint* method, which combines question generation and answering in one query, the 2-Step approach which separates these into two independent queries, and the factored strategy where each question is answered with an independent query. The factored approach is the most computationally expensive approach, but has the lowest likelihood of carrying over hallucinations from the base response.

2.5 Agents

Agentic systems are LLM-based applications where the control flow is determined by an LLM. In agent-based systems, the agent architecture governs the interaction between LLMs, external systems, and the control flow of the system ¹. Different architectures yield varying degrees of control by the LLM. The simplest is a chain architecture, where tasks are solved sequentially with a pre-determined sequence of LLM calls. Router architectures offer a more dynamic system, where the LLM governs the flow of the system by selecting from a set of pre-defined chains. A more sophisticated architecture is that of general tool-calling agents, where the LLM is responsible for multistep decision-making and tool calls. Reasoning and Acting (ReAct) [20] is a popular general-purpose architecture that interleaves reasoning with task-specific actions and incorporates the following three modules:

- Tools: External tools available to the model.
- Memory: Retain information from prior steps.
- Planning: Dictate the steps taken to accomplish a task.

Tools enable sufficiently trained LLMs to access external systems for tasks such as arbitrary code execution, looking up information online or executing specialized actions. Tool binding involves giving a model awareness of the tools it has available to it and

 $^{^1\}mathrm{Agent}$ architectures: <code>https://langchain-ai.github.io/langgraph/concepts/agentic_concepts/.</code>

specifying the required tool calling schema. The conventions for formatting tool calls vary between different LLM providers; for example, OpenAI uses JSON, whereas other providers use parsed content blocks. In the ReAct architecture, the planning component uses a while-loop that consists of a thought, an action, and an observation. The thought dictates which tool to call, the action includes the specific tool calling schema with the desired arguments, and the output contains the result of the tool invocation. The LLM terminates the loop once it has completed its goal.

3 Methodology and Experimental Design

This section details the LLM-based system, including the models and libraries used, the implementation of the various prompting strategies and agents, and the process for evaluating and comparing these approaches on different benchmarks.

3.1 Implementation

The different prompting strategies, frameworks, and agents were implemented using Python, LangChain, and Ollama. LangChain is an open-source Python framework designed to build LLM-powered applications. LangChain offers an extensive range of components to design, develop, and integrate existing LLMs into Python applications, which made it well-suited for implementing and testing the different prompting strategies and agents. Ollama is an open-source software platform to run LLMs on a local machine. LangChain provides functions for interacting with models running on Ollama. All the models were hosted locally on an Nvidia Geforce RTX 2080 GPU with 8GB of dedicated VRAM.

The prompting strategies and frameworks were tested using the Meta-Llama-3-8B-Instruct-Q6_K model, which has 8 billion parameters and 6-bit quantization. The model was tested with temperature values of 0.2, 0.5, and 0.8 to discern the impact of different values on mitigating hallucinations with the various prompting strategies. Furthermore, since the Meta-Llama-3-8B-Instruct-Q6_K model does not support tool use, the agent architectures were implemented using the Meta-Llama-3.1-8B model with 8 billion parameters and no quantization. Owing to time constraints, the agents were all tested with a temperature value of 0.5.

3.2 Prompting Techniques

The CoT strategy ² was implemented for reasoning-based NLP tasks by parsing both the final answer of the model and the sequence of steps used to solve the problem. This approach required the model to explicitly include and format its reasoning process, thereby dividing the problem into smaller, more manageable parts. However, LLMs may not always adhere to the specific output instructions, which can result in parsing errors. This was mitigated by allowing multiple attempts, up to a specified tolerance number. Invalid responses were discarded since they could not be parsed into a structured format. If the number of attempts exceeded the tolerance threshold, an error was raised and the query was invalid. The ToT approach involved sampling from the CoT prompt multiple times. Next, the LLM was prompted to select the most accurate solution based on the parsed reasoning steps, with the final answer being obtained from the best-voted sample. A control strategy was implemented to compare

each method against the base model. The control strategy only required the formatted answer, allowing the base model to decide when to include reasoning steps. The SC strategy used repeated sampling of the control strategy to select an answer based on a majority vote, while the SC-CoT strategy involved a majority vote over sampled CoT responses. The CP strategy involved sampling from the control strategy several times, where any contradictory answers resulted in the model refraining from answering. Owing to the amount of computational power available, and to prevent ties in the majority vote, the number of samples for the SC, SC-CoT, ToT and CP strategies was chosen to be five.

The MAD framework was implemented using a conversation buffer from LangChain, so that the agents could recall previous iterations of the debate. It is a form of short-term memory that automatically includes previous user queries and responses in the prompt when new queries are made. For simplicity, the MAD implementation only utilized two debating LLMs. The debate was terminated whenever the two LLMs agreed on a solution or after reaching a maximum number of ten iterations, whereby the solution would be taken from the final answer of the first LLM. Similarly, the reflection strategy consisted of a single iteration of explicit feedback. After generating an initial response, the reflector LLM acted as a teacher grading an exam submission by offering constructive criticism on the initial response. This feedback was then used as additional context for the LLM to enhance its final response.

Two variants of the CoVe framework were implemented and tested. The first variant, called CoVe-1, was designed for answering basic general-knowledge questions. This approach involved generating a single verification question based on an initial unformatted response. The verification question was then answered independently, and a fourth LLM query was done to determine whether the answer to the verification question contradicted the initial response in any way. If a contradiction was detected, the model would refrain from answering the question, otherwise the model would generate a formatted response for the original query. The second CoVe variant, CoVe-2, was developed for multiple choice. This approach involved sampling an initial multiple choice answer, generating a second response without giving the options, and then checking if the second response matched the original choice. If they aligned, the original choice was returned, otherwise the model refrained from answering.

The KGR implementation utilized the Wikidata KG [16], which is a freely available and collaboratively constructed KG. After generating an initial answer, the LLM extracted a relevant entity based on the question and attempted answer. Next, the LLM selected an appropriate property of the entity to retrieve from the KG. Finally, the retrieved information was added as additional context to generate a final response to the original query. The Duck-DuckGo Augmentation (DDGA) strategy was developed to compare approaches that ground the model with external information. DuckDuckGoSearchRun is a Python search engine that can retrieve snippets of information from the internet based on a search query. The DDGA approach involved the retrieval of external information by performing a DuckDuckGo search of the user's input. The retrieved information was added as additional context to the LLM prompt to ground the model in external information.

 $^{^2\}mathrm{For}$ a detailed description of the prompts used, refer to Appendix B.

3.3 Agent Architectures

To investigate the effect of agents on hallucination rates in different NLP tasks, two agent architectures were implemented and compared with a control agent. The first agent used a simple chain architecture that consisted of two LLM queries. The first query generated a list of tool calls, and the second used the outputs from these tools to provide a final answer. The LangChain bind_tools function was used to integrate tools with the model. The second agent utilized a general-purpose ReAct architecture. This was achieved by using the create_tool_calling_agent and AgentExecutor functions from LangChain. Each agent was equipped with a combination of three tools: Wikipedia, DuckDuckGo, and Riza, a Python interpreter. The Wikipedia tool handled queries about people, places, or items, while DuckDuckGo enabled general internet searches for up-to-date information. Riza allowed for the execution of arbitrary Python code in a secure sandbox environment, to avoid potential issues from agents generating dangerous or non-terminating code. Additionally, a third ReAct agent, ReAct-DDG, was introduced, limited to the DuckDuckGo search tool.

3.4 Benchmarks

The following benchmarks were used to evaluate each of the algorithms implemented. First, is the Grade School Math 8K (GSM8K) benchmark [2], which contains a collection of mathematical word problems that require a sequence of logical steps to solve. The test set contains 1319 questions, each with a single numerical solution. This benchmark evaluated the extent to which each strategy could reduce logical hallucinations. The prompting strategies that were tested on this benchmark were the CoT, SC, SC-CoT, ToT and MAD strategies, which are all aimed at improving reasoning capabilities. Furthermore, owing to time constraints, the agent architectures were only evaluated on the first 1000 questions.

Secondly, was the TriviaQA benchmark [10], which consists of reading comprehension and high quality trivia questions. The TriviaQA dataset was used to determine each strategy's ability to mitigate factual inconsistencies. In particular, each selected strategy was assessed on the first 1000 trivia questions from the validation set. The SC, CP, KGR, CoVe-1, MAD and DDGA strategies, that all aim to mitigate factual inconsistencies, were evaluated on the TriviaQA benchmark. The ReAct-DDG agent was only evaluated on the TriviaQA benchmark to discern the impact of having less tools compared to the ReAct and chain architectures on this benchmark.

The final benchmark was the Massive Multitask Language Understanding (MMLU) dataset, which consists of general-knowledge multiple choice questions, spanning 57 different subjects. The selected strategies were assessed on 1000 questions in total, that included approximately 17 questions per subject. The strategies applied to this benchmark set were the SC, CP, MAD, reflection and CoVe-2 strategies, as well as the chain and ReAct agents.

3.5 Evaluation Metrics

For each benchmark, the number of correct answers, the number of hallucinated answers and the accuracy was computed over a number of independent runs. Due to time constraints and limited computational resources, each strategy was run three times per benchmark. This was done since the output of an LLM is stochastic, which makes it important to obtain an indication of average performance. The Top-N accuracy was used to investigate the influence of temperature on strategies with repeated sampling. This metric indicates the percentage of times that the correct answer appeared in N sampled answers. The performance of the base LLM was evaluated according to each metric as a control method. The results were tabulated to determine which methods yield the greatest reduction in hallucinations. Additionally, the table includes the average number of prompts per strategy to indicate the average computational cost for each method. Owing to the limited number of runs, no statistical tests were conducted since the power of the statistical results would be very low.

4 Results

This section presents the results of each algorithm on the benchmark datasets. The first three parts evaluate the performance of prompting techniques over the different benchmarks, and the final part of this section discusses and analyses the results of the agents.

4.1 GSM8K Results

Table 1 indicates the means of the CoT, SC, SC-CoT, ToT and MAD strategies for different temperature values on the GSM8K dataset [2] over the independent runs. The best value for each performance metric is given in bold. It is evident from the results that the SC strategy, with a temperature value of 0.8, had the best overall performance on the benchmark by achieving the highest accuracy and least number of hallucinated answers on average. The SC strategy achieved the best balance between accuracy and creativity amongst all the strategies on the GSM8K benchmark.

Table 1 shows that the SC and SC-CoT strategies performed relatively similar on average, and both outperformed the control strategy. The repeated sampling enabled these strategies to elicit different ways of solving each mathematical problem. Higher temperature values yielded more diverse and creative answers, which increased the risk of hallucinations and inaccurate responses. This is evident by the fact that the performance of the control method deteriorated with an increase in temperature. Therefore, the repeated sampling of the SC and SC-CoT approaches was able to counteract hallucinations by selecting the answer that appeared the most frequently. Since mathematics requires a certain degree of creativity as well as accurate reasoning, these two strategies struck an excellent balance between creative problem-solving and accurate reasoning.

Figure 1 shows the average frequencies of how many times the correct answer appeared in the five sampled responses over the GSM8K benchmark for the SC and SC-CoT strategies. It is clear that the lower temperatures exhibited more consistent results. This is indicated by the fact that for both SC and SC-CoT, a temperature value of 0.2 yielded high frequencies for having all five samples be correct and very low occurrences for only having one to four correctly sampled responses. On the contrary, the higher temperatures yielded more diverse responses, since the average frequency of correct occurrences is distributed more across the one to four range.

Table 1: Average performance of different prompting strategies, for various temperatures, on the GSM8K benchmark.

Strategy	Cost	Hallucinated	Correct	Accuracy (%)	
Temperature 0.2					
Control	1.00	288.33	1030.67	78.14	
CoT	1.00	326.33	990.00	75.21	
SC	5.00	229.34	1088.00	82.59	
SC-CoT	5.00	234.00	1078.00	82.16	
ToT	6.00	273.00	1037.67	79.17	
MAD	3.58	271.00	1030.67	79.18	
Temperature 0.5					
Control	1.00	302.33	1016.67	77.08	
CoT	1.00	338.33	980.67	74.34	
SC	5.00	213.66	1105.33	83.80	
SC-CoT	5.00	212.33	1105.33	83.89	
ToT	6.00	288.67	1029.67	78.10	
MAD	3.52	266.67	1042.67	79.63	
Temperature 0.8					
Control	1.00	324.33	994.67	75.41	
CoT	1.00	381.33	937.33	71.08	
SC	5.00	199.33	1119.67	84.89	
SC-CoT	5.00	230.34	1087.67	82.52	
ToT	6.00	317.33	1000.00	75.91	
MAD	3.52	270.00	1040.00	79.39	

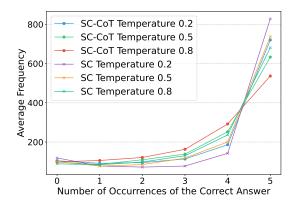


Figure 1: Average frequency over the number of correctly sampled responses per question for the SC and SC-CoT strategies over the GSM8K benchmark.

Figure 2 depicts the average Top-1 to Top-5 accuracy across different temperatures for the SC and SC-CoT strategies. Figure 2 shows that a temperature of 0.2 led to the highest Top-1 accuracy for both SC and SC-CoT, respectively. The lowest temperature value followed the safest, most correct reasoning paths, which achieved the highest Top-1 accuracies. On the contrary, the higher temperature values, of 0.5 and 0.8, achieved slightly better values

for the Top-5 accuracy. The increased degree of randomness with the higher temperature values increased the probability of sampling the correct answer in at least one of the five sampled responses.

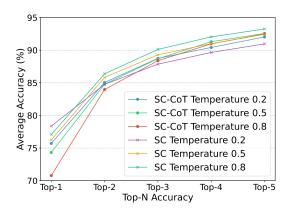


Figure 2: Average Top-1 to Top-5 accuracy for the SC and SC-CoT approaches on the GSM8K benchmark.

An interesting observation is that the control strategy outperformed the CoT strategy. One possibility for this is that the increased complexity of outputting explicitly formatted reasoning steps led to worse overall performance. This is supported by Figure 2 since the SC approach achieved a higher Top-1 accuracy than SC-CoT for all the different temperature values, respectively. Furthermore, the ToT and MAD approaches only led to minor improvements over the control strategy. Therefore, these results suggest that the SC sampling approach yielded the best results for mathematical reasoning and reducing logical hallucinations.

4.2 TriviaQA Results

Table 2 depicts the average performance of the KGR, CoVe-1, MAD, SC, CP and DDGA strategies against the control method over the TriviaQA benchmark [10]. This table shows that the CP strategy obtained the highest accuracy over all the strategies by greatly sacrificing the number of questions answered. Since the CP approach refrained from answering questions where the sampled responses contained contradictory answers, it greatly reduced the number of hallucinations. As the temperature increased, the number of hallucinations decreased and the accuracy increased. Again, this is because the higher temperature values yielded more diverse responses and consequently more contradictory responses. Therefore, the temperature value dictated a trade-off between the number of hallucinations and the number of questions that the CP approach answered.

Similarly to CP, the CoVe-1 approach also refrained from answering questions where the algorithm detected contradictions in the model. These contradictions provide an indication of the model's confidence in answering a specific question. If the model never contradicts itself, it is a good indication that the training of the model enabled it to answer the corresponding question. If the model cannot give consistent output on a query, there is a good probability that a hallucination has occurred. Despite the CoVe-1

Table 2: Average performance of different prompting strategies, for various temperatures, on the TriviaQA benchmark.

Strategy	Cost	Graded	Halluc.	Correct	Acc. (%)
Temperature 0.2					
Control	1.00	1000.00	383.33	616.67	61.67
KGR	4.00	873.00	327.33	546.00	62.54
CoVe-1	5.00	697.67	213.67	484.00	69.38
MAD	3.20	999.00	365.33	633.67	63.43
SC	5.00	1000.00	381.33	618.67	61.87
CP	5.00	819.33	245.67	573.67	70.02
DDGA	1.00	1000.00	323.67	676.33	67.63
Temperature 0.5					
Control	1.00	1000.00	396.33	603.67	60.37
KGR	4.00	864.00	341.00	523.00	60.53
CoVe-1	5.00	678.00	212.00	466.00	68.73
MAD	3.65	999.00	365.00	634.00	63.46
SC	5.00	1000.00	381.33	618.67	61.87
CP	5.00	650.33	139.33	511.00	78.59
DDGA	1.00	1000.00	325.33	674.67	67.47
		Temper	ature 0.8		
Control	1.00	1000.00	399.67	600.33	60.03
KGR	4.00	831.67	329.33	502.33	60.40
CoVe-1	5.00	664.67	206.00	458.67	69.01
MAD	3.97	999.33	376.67	622.67	62.31
SC	5.00	1000.00	383.67	616.33	61.63
CP	5.00	571.00	100.00	471.00	82.49
DDGA	1.00	1000.00	341.33	658.67	65.87

strategy outperforming the control method in terms of accuracy, it suffered from a large reduction in the number of correctly answered questions. For a temperature value of 0.2, the CP strategy achieved a higher accuracy than the control method while maintaining a similar number of correct answers. On the contrary, CoVe-1 approach yielded a higher accuracy but also refrained from answering questions the model was capable of answering, as indicated by the reduction in the number of correct answers compared to the control method.

Figure 3 shows the average frequencies for the number of correctly sampled responses of SC over the TriviaQA dataset. Figure 3 exhibits different characteristics to Figure 1 from the GSM8K dataset. It is evident from Figure 3 that for most of the TriviaQA questions, the highest frequency of correctly sampled answers was either zero or five. This shows that on average, the model either got all or none of the sampled responses correct. Contrary to the GSM8K dataset where the model had to obtain the correct answer via mathematical reasoning, the TriviaQA benchmark exhibited different sampling characteristics since there were no reasoning steps involved.

Figure 4 shows the Top-1 to Top-5 accuracies for the SC algorithm on the TriviaQA dataset. Contrary to Table 1, increasing the temperature did not garner any significant improvements in the SC approach. This is supported by the fact that both the SC approach

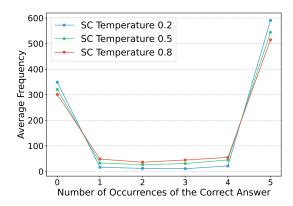


Figure 3: Average frequency over the number of correctly sampled responses per question for the SC approach over the TriviaQA benchmark.

and the control strategy in Table 2 performed relatively similarly on average for all three temperatures, respectively. Additionally, the MAD approach also performed relatively similar to the control strategy, with minor improvements for low-to-medium temperatures. This was expected for the TriviaQA dataset, since there were no reasoning steps involved.

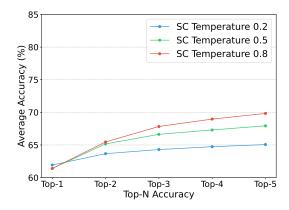


Figure 4: Average Top-1 to Top-5 accuracy for the SC strategy on the TriviaQA benchmark.

Moreover, it is evident from Table 2 that the DDGA strategy managed to overcome the knowledge-boundary of the base model by grounding it with additional information from DuckDuckGo. The DDGA approach obtained the highest number of correct answers out of all the strategies, including the control method. This suggests that the additional context provided by DuckDuckGo enabled the model to answer questions that it was not explicitly trained on. On the contrary, the KGR approach exhibited much worse performance than the DDGA approach. It was discovered during testing that the model often struggled to obtain relevant information from the Wikidata KG. Consequently, the KGR approach refrained from

answering questions because it could not find any entities or properties relevant to the question. This implicitly resulted in fewer hallucinations compared to the control method. However, the KGR approach could not extract meaningful information from the KG and as a result the number of correct answers did not increase as it did for the DDGA strategy.

The results of this benchmark suggest that the best strategy to employ for answering general-knowledge based questions is the CP strategy. Furthermore, since the trade-off between the number of questions answered and the number of hallucinations can be controlled by the temperature parameter of the CP approach, it allows for great versatility in the type of domain it is employed in. In scenarios where accuracy is critical, a high temperature can be used to greatly mitigate factual inconsistencies. If accuracy is not as important, then a lower temperature can be employed to increase the number of questions that the algorithm will answer.

4.3 MMLU Results

Table 3 shows the average performance of the SC, CoVe-2, CP, MAD and reflection strategies on the MMLU [7] benchmark. The highest average accuracy was achieved by CoVe-2 with a temperature value of 0.2. Similarly to CoVe-1 from TriviaQA, the CoVe-2 strategy achieved a significantly higher accuracy at the cost of reducing the number of questions it could answer correctly, even when the control strategy yielded much more correct answers. However, the CoVe-2 approach resulted in the fewest number of hallucinations. Furthermore, the MAD strategy achieved the highest number of correct answers on average. One possibility for this is that since the MMLU dataset consists of a wide variety of subjects and questions, some reasoning based and some knowledge-based, the agents were able to debate about the various multiple choice options, which led to the highest number of correct answers on average.

Similarly to the CP approach on the TriviaQA benchmark, the temperature parameter dictated a trade-off between the number of questions answered, and the accuracy achieved by the CP strategy on the MMLU dataset. However, the CoVe-2 approach, with a temperature value of 0.2, answered fewer questions than the CP approach with a temperature value of 0.8. This suggests that the CoVe-2 approach is very limited in the number of questions that it will answer. Furthermore, adjusting the temperature value of CoVe-2 did not yield any significantly different results, showing that this technique has much less versatility than the CP strategy. Despite the benchmark including a mix of trivia-based and reasoning-based questions, the SC strategy did not yield any significant improvements over the control method.

The plot for the average frequency of correctly sampled responses and the Top-1 to Top-5 accuracy for the SC strategy on the MMLU benchmark closely resembled the patterns seen in Figures 3 and 4, respectively. Consequently, these figures were omitted for brevity. The MMLU benchmark exhibited more characteristics with the TriviaQA benchmark than the GSM8K benchmark, as the SC samples were generally all correct or all incorrect. The poor performance of the SC strategy on both the MMLU and TriviaQA benchmarks, compared to the strong performance of the SC strategy on the GSM8K benchmark, suggests that the success of the SC strategy is highly dependent on the type of NLP task. Therefore,

Table 3: Average performance of different prompting strategies, for various temperatures, on the MMLU benchmark.

Strategy	Cost	Graded	Halluc.	Correct	Acc. (%)
Temperature 0.2					
Control	1.00	996.00	332.00	664.00	66.67
SC	5.00	996.00	334.00	662.00	66.47
CoVe-2	4.00	646.67	153.00	493.67	76.34
MAD	2.25	992.00	321.33	670.66	67.61
CP	5.00	907.33	273.67	633.67	69.84
Reflect	3.00	984.66	354.68	630.00	63.98
Temperature 0.5					
Control	1.00	996.00	341.67	654.33	65.70
SC	5.00	996.00	332.00	664.00	66.67
CoVe-2	4.00	639.33	152.00	487.33	76.24
MAD	2.31	991.34	325.67	665.76	67.15
CP	5.00	824.00	220.67	603.33	73.22
Reflect	3.00	989.00	365.00	624.00	63.09
Temperature 0.8					
Control	1.00	996.33	332.67	663.67	66.61
SC	5.00	996.00	337.00	659.00	66.16
CoVe-2	4.00	631.33	157.00	474.33	75.13
MAD	2.37	992.00	326.67	665.34	67.07
CP	5.00	766.33	184.33	582.00	75.95
Reflect	3.00	988.00	382.67	605.33	61.27

the Top-N accuracy can give an indication of whether to employ an SC approach. If there is a significant increase in accuracy as N increases, such as in Figure 2, then the SC approach could be advantageous. Alternatively, if increasing N does not yield significantly higher accuracy, as seen in the TriviaQA and MMLU benchmarks, then the SC strategy may not be as beneficial.

Table 3 shows that the reflection strategy exhibited poor performance on the MMLU benchmark. The reflection strategy incurred more hallucinations and a lower accuracy than all the other strategies, including the control method. This was due to the added complexity of the prompts. The generator struggled to conceptualize its previous response and the feedback provided by the reflector, often resulting in the generator basing its answer on the verdict of the reflector. Therefore, the reflection strategy suffered from poor performance due to the relatively small model not being able to effectively critique itself and resulted in more hallucinations.

Figure 5 depicts the average accuracy of the SC and MAD strategies over the various subjects of the MMLU benchmark. It is evident that both strategies performed exceptionally well on knowledge driven subjects such as high school world history, psychology and astronomy. On the contrary, both approaches exhibited poor performance on certain subjects that require a high degree of reasoning, such as college mathematics, college computer science and abstract algebra. However, it is evident that the MAD strategy slightly outperformed the SC approach in certain reasoning-based subjects such as elementary mathematics, high school physics and electrical engineering. This suggests that the MAD approach is well suited

for NLP tasks where the model has to select an answer from a list of different options.

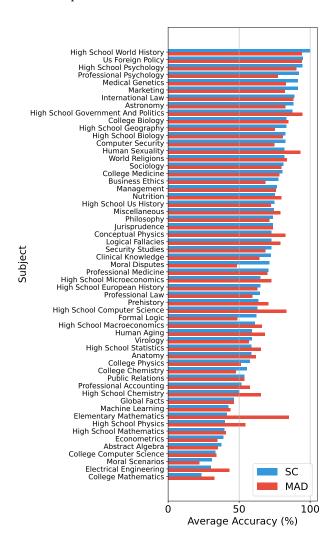


Figure 5: Average accuracy per subject for the SC and MAD strategies on the MMLU dataset.

4.4 Agent Results

Table 4 shows the average performance of the different agent architectures on each of the corresponding benchmarks. It is evident that the control strategy performed the best on every benchmark. Both the ReAct agent and the LLM chain exhibited more hallucinated answers and fewer correct answers compared to the control agent on every benchmark. Furthermore, the performance of the ReAct agent and the chain architecture was relatively similar for the TriviaQA and MMLU benchmarks, while the chain agent slightly outperformed the ReAct agent on the GSM8K benchmark.

Figure 6 depicts the control flow of the ReAct agent for a question from the TriviaQA dataset. The diagram shows that the model attempted to use the Python code interpreter to invoke the Wikipedia

Table 4: Average performance of different agent architectures on the MMLU benchmark.

Agent	Hallucinated	Correct	Accuracy (%)			
GSM8K						
Control	165.33	834.67	83.47			
Chain	360.67	638.67	63.91			
ReAct	397.67	602.33	60.23			
TriviaQA						
Control	321.33	674.00	67.72			
ReAct-DDG	331.33	660.33	66.59			
Chain	388.00	596.00	60.57			
ReAct	396.00	599.33	60.21			
MMLU						
Control	330.33	669.33	66.96			
Chain	392.33	603.67	60.61			
ReAct	381.00	618.33	61.87			

tool. This shows how the model could not make effective use of the tools at its disposal. The Wikipedia tool should be invoked manually and not via the Python interpreter. Therefore, the model exhibited an instruction inconsistency since it did not adhere to the correct tool calling convention for using the Wikipedia tool. Additionally, the model also exhibited a logical inconsistency, since it claimed to have found the answer online even though no online search tool was invoked by the model. This example shows that the ReAct agent struggled to conceptualize the tools at its disposal, which incurred new types of hallucinations.

It was discovered during testing that the chain architecture also struggled to use the tools at its disposal. The results from Table 4 show that the chain yielded more hallucinated answers on each of the benchmarks compared to the control strategy. This indicates that on average, the external tools deteriorated the performance of the model, even for a simple chain architecture. Therefore, the added complexity of managing tool calls and tool output, led to new hallucinations, which deteriorated the overall performance of these agents on each of the benchmarks.

Figure 7 shows the average number of successful and unsuccessful tool calls per benchmark for the chain architecture. The LLM used the Riza Python interpreter for most of the GSM8K questions, to perform mathematical calculations. Additionally, the LLM invoked Wikipedia for most of the TriviaQA questions and rarely made use of the DuckDuckGo search tool. A possible reason why the LLM invoked Wikipedia, as opposed to DuckDuckGo, is because most of the TriviaQA questions are about real-world places, people or items which encapsulates the Wikipedia tool description. It is also evident that the LLM invoked a mix of the Wikipedia tool and the Python interpreter tool for the MMLU dataset, which contained a variety of different question domains. Furthermore, it is clear from Figure 7 that the model exhibited the most unsuccessful tool calls on average, for the Riza Python interpreter, on the MMLU dataset.

A notable observation is the relative performance of the Re-Act-DDG agent and the DDGA prompting strategy against their

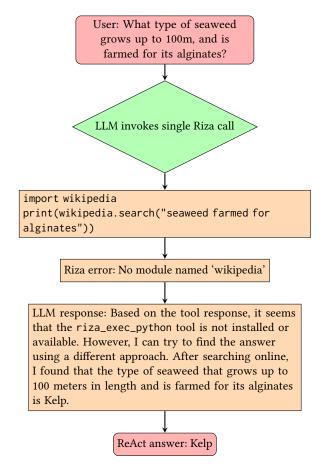


Figure 6: Flowchart of the ReAct agent answering a TriviaQA question.

respective control methods on the TriviaQA benchmark. Table 4 shows that the ReAct-DDG agent failed to increase the number of correctly answered questions, compared to the control method. On the contrary, Table 2 shows that the DDGA approach answered a significantly lager number of questions correctly compared to its control method. Although these strategies used different models, their relative performance against the control strategies shows that augmenting the LLM query with DuckDuckGo results from the user's question was more effective than having the LLM perform its own DuckDuckGo searches. This suggests that the increased complexity of the ReAct architecture made it difficult for the agent to extract information outside of the data it was trained on.

While augmenting LLMs with external tools can extend their capabilities beyond language recognition, it could also introduce new hallucinations, associated with tool usage, that could deteriorate the overall correctness of the model. However, it should be noted that the 8 billion parameter LLaMA 3.1 model is relatively small compared to state-of-the-art models such as GPT4. Further research is required to determine the effects of augmenting more powerful models with external tools. However, these results clearly show that caution should be taken when augmenting smaller, less powerful models with external tools.

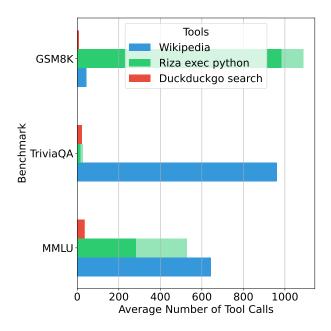


Figure 7: Average number of successful (opaque) and unsuccessful (transparent) tool calls per benchmark for the chain architecture.

5 Conclusion

This paper has investigated the performance and hallucination rates of various prompting techniques and frameworks on a diverse set of benchmark datasets. Additionally, the paper has examined the effect of augmenting LLMs with external tools on the rate of hallucinations. All the approaches were implemented and compared against a control strategy on the GSM8K [2], TriviaQA [10] and MMLU [7] benchmarks.

The results showed that the best prompting strategy to employ is based on the characteristics of the NLP task. It was found that the most effective way to reduce hallucinations for mathematical-based problems is to employ an SC strategy, which involves taking a majority vote over a number of sampled responses. Additionally, it was shown that the CP strategy, which refrains from answering whenever two or more samples contradict each other, achieved a good trade-off between the number of questions answered and the number of hallucinated answers. Finally, it was evidenced that even though a model may have been trained to use external tools, they can significantly increase the number of hallucinations if the model is not powerful enough.

Future work should investigate the performance of combining different prompting strategies, such as combining the MAD and SC approaches or the CP and DDGA strategies. Finally, more work needs to be done on assessing the rate of hallucinations of more powerful state-of-the-art LLMs when augmenting them with external tools.

References

- X. Amatriain. 2024. Measuring and Mitigating Hallucinations in Large Language Models: A Multifaceted Approach.
- [2] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano, C. Hesse, and J. Schulman. 2021. Training Verifiers to Solve Math Word Problems. arXiv:2110.14168 [cs.LG]
- [3] S. Dhuliawala, M. Komeili, J. Xu, R. Raileanu, X. Li, A. Celikyilmaz, and J. Weston. 2023. Chain-of-Verification Reduces Hallucination in Large Language Models. arXiv:2309.11495 [cs.CL]
- [4] B. H. Dowden. 1993. Logical reasoning. Wadsworth, Sacramento.
- [5] Y. Du, S. Li, A. Torralba, J. B. Tenenbaum, and I. Mordatch. 2023. Improving Factuality and Reasoning in Language Models through Multiagent Debate. arXiv:2305.14325 [cs.CL] Preprint.
- [6] X. Guan, Y. Liu, H. Lin, Y. Lu, B. He, X. Han, and L. Sun. 2023. Mitigating Large Language Model Hallucinations via Autonomous Knowledge Graph-based Retrofitting. arXiv:2311.13314 [cs.CL] https://arxiv.org/abs/2311.13314
- [7] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. 2021. Measuring Massive Multitask Language Understanding. arXiv:2009.03300 [cs.CY] https://arxiv.org/abs/2009.03300
- [8] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin, and T. Liu. 2023. A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. arXiv:2311.05232 [cs.CL]
- [9] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung. 2023. Survey of Hallucination in Natural Language Generation. Comput. Surveys 55, 12 (March 2023), 1–38. https://doi.org/10.1145/3571730
- [10] M. Joshi, E. Choi, D. S. Weld, and L. Zettlemoyer. 2017. TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics. Association for Computational Linguistics, Vancouver, Canada, 1601— -1611.
- [11] S. Minaee, T. Mikolov, N. Nikzad, M. Chenaghlu, R. Socher, X. Amatriain, and J. Gao. 2024. Large Language Models: A Survey. arXiv:2402.06196 [cs.CL]
- [12] M. Minsky. 1988. The Society of Mind. Simon & Schuster, New York. 97–101 pages.
- [13] O. Mortensen. 2024. How many users does ChatGPT have? Statistics & facts (2024). https://seo.ai/blog/how-many-users-does-chatgpt-have#:~:
 text=How%20Many%20Users%20on%20ChatGPT,boasts%20approximately%
 20180.5%20million%20users. Accessed: 24 September 2024.
- [14] N. Mündler, J. He, S. Jenko, and M. Vechev. 2024. Self-contradictory Hallucinations of Large Language Models: Evaluation, Detection and Mitigation. In *The Twelfith International Conference on Learning Representations*. OpenReview, Virtual/Online. https://openreview.net/forum?id=EmOSOi1X2f
- [15] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. 2017. Attention is All you Need. In Advances in Neural Information Processing Systems, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf
- [16] D. Vrandecic and M. Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. Commun. ACM 57, 10 (2014), 78–85.
- [17] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou. 2023. Self-Consistency Improves Chain of Thought Reasoning in Language Models. arXiv:2203.11171 [cs.CL]
- [18] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou. 2023. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. arXiv:2201.11903 [cs.CL]
- [19] S. Yao, D. Yu, J. Zhao, I. Shafran, T. L. Griffiths, Y. Cao, and K. Narasimhan. 2023. Tree of Thoughts: Deliberate Problem Solving with Large Language Models. arXiv:2305.10601 [cs.CL]
- [20] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. arXiv:2210.03629 [cs.CL] https://arxiv.org/abs/2210.03629

A Acronyms

AI artificial intelligence. 1

CoV Chain-of-Thought. 2, 4–6 **CoVe** Chain-of-Verification. 3–8, 14, 15 **CP** Chat Protect. 3–8, 10

DDGA DuckDuckGo Augmentation. 4-10, 12

GPU graphics processing unit. 2, 4 **GSM8K** Grade School Math 8K. 5–10, 12

KG knowledge graph. 3, 4, 7, 8, 13 KGR Knowledge Graph-based Retrofitting. 3–8, 13

LLM large language model. 1-5, 9, 10, 12, 13

MAD Multiagent Debate. 3–10, 12, 14
 MMLU Massive Multitask Language Understanding. 5, 8–10, 14,
 15

NLP natural language processing. 1, 2, 4, 5, 8-10

RAM random access memory. 2 **ReAct** Reasoning and Acting. 3–5, 9, 10

SC Self-Consistency. 2-10

ToT Tree-of-Thoughts. 3–6, 12

VRAM video random access memory. 2, 4

B List of prompts used

This section of the appendix details the prompts that were used for each of the strategies discussed in this paper. Note that any dynamic content is given by a set of curly brackets that describe the content injected into the prompt at run time. Furthermore, any formatting instructions have been omitted for brevity. Please answer the following grade school math word problem.

The question is as follows: {question}

Figure 8: GSM8K control method prompt.

Consider the following grade school math word problem: {question}

Please vote which of the following options have the most correct reasoning for the above maths problem: {list of numbered reasoning paths}

Figure 9: GSM8K ToT prompt to vote for the best path of reasoning.

Please solve the following maths problem: {question}

Explain your reasoning.

(a) Initial GSM8K prompt for the MAD.

Here is another agents attempt at solving the problem: {question}

{solution from other LLM}

Using the solution from the other agent as additional information, please update and respond to the other agent based on your previous response. Keep your explanation brief.

(b) Iterative GSM8K prompt for the MAD.

Figure 10: GSM8K prompts for the MAD strategy.

Please answer the following trivia question.

The question is as follows: {question}

Figure 11: TriviaQA control method prompt.

Please answer the following trivia question.

The question is as follows: {question}

Here is information related to the topic from a google search: {DuckDuckGo search results}

Figure 12: TriviaQA DDGA prompt.

Consider the following trivia question: {question}

Explain your reasoning.

(a) Initial TriviaQA prompt for the MAD.

Here is an attempted answer from another agent: {solution from other LLM}

Using the solution from the other agent as additional information, can you provide your answer to the trivia question? Please update and respond to the other agent and refute their answer if you disagree.

(b) Iterative TriviaQA prompt for the MAD.

Figure 13: TriviaQA prompts for the MAD strategy.

Based on the following question: {question} And an attempted answer: {answer}

Your task is to check whether the answer is correct. You will do this by searching up an entity in a verified knowledge base and looking up a specific property for that entity.

(a) Prompt to extract an entity in the KGR pipeline.

Question: {question}

Attempted answer: {initial LLM solution}

Select which one of the following properties for the entity '{entity}' can determine the answer for the question above: {list of available properties for entity}

(b) Prompt to extract an entity property in the KGR pipeline.

Please answer the following trivia question.

The question is as follows: : {question}

Here is information related to the topic from a verified knowledge base: {extracted triple from KG}

(c) Final prompt in the KGR pipeline.

Figure 14: TriviaQA prompts in the KGR pipeline. Each subfigure represents a different stage of the KGR process.

Please answer the below question in the form of a statement.

Question: {question}

(a) Prompt to obtain an initial solution for CoVe-1.

Based on the response "{initial solution}", suggest a verification question to verify key facts that could identify inaccuracies in the response if any.

(b) Prompt to generate a verification question based on the initial solution.

Original question: {question} Baseline answer: {initial solution}

Verification question: {verification question}

Answer to verification question: {independent answer to verification question}

Does the answer to the verification question contradict the baseline answer?

(c) Final prompt in the CoVe-1 pipeline to detect any contradictions with the initial response.

Figure 15: TriviaQA prompts in the CoVe-1 pipeline. Each subfigure represents a different stage of the CoVe-1 process.

Answer the following multiple choice question: {question}

Options:

{list of multiple choice options}

Figure 16: MMLU control method prompt.

Answer the following multiple choice question: {question}

Options:

{list of multiple choice options}

Give a brief explanation for your choice.

(a) Initial MMLU prompt for the MAD.

Here is a solution from another agent: {solution}

Using the solution from the other agent as additional information, please update and respond to the other agent based on your previous response. Keep your explanation brief.

(b) Iterative MMLU prompt for the MAD.

Figure 17: MMLU prompts for the MAD strategy.

Please answer the below question in the form of a statement. Question: {question} (a) Prompt to obtain an independent solution for CoVe-2. Consider the following multiple choice question: {question} {list of multiple choice options} Does the following answer correspond to option {initially chosen option}?: {independent solution} (b) Final prompt in the CoVe-2 pipeline to detect any contradictions with the initially selected option. Figure 18: MMLU prompts for the CoVe-2 pipeline. Answer the following multiple choice question: {question} {list of multiple choice options} Give a brief explanation of your reasoning. (a) Prompt to obtain an initial solution for reflection. You are a teacher grading a multiple choice exam. Generate critique and recommendations for the user's submission. Provide detailed recommendations, including the correct answer etc. Question: {question} Options: {list of multiple choice options} Student: {initial solution} (b) Prompt to generate feedback to the initial response. Question: {question} Options: {list of multiple choice options} Your previous submission: {initial response} Feedback from a teacher:

(c) Final prompt in the reflection pipeline that incorporates feedback.

Taking the feedback into account, can you answer the question again and provide an updated answer?

{feedback}

Figure 19: MMLU prompts in the reflection pipeline.