

Influence of Artificial Degrees of Certainty in Large Language Model Reasoning-Based Tasks

Background

- Large language models (LLMs) are artificial intelligence models that use machine learning to understand and generate human language
 - Ex. OpenAI ChatGPT, Google Gemini, Microsoft Copilot
- As LLMs continue to be integrated into decision-making systems (e.g. medical diagnosis, customer service), it becomes imperative to understand the factors that influence their performance
- Prior research has demonstrated that various human inputs—such as specific details (intuitively), subtle biases, or even tone of voice—can significantly shape LLM responses

Prompt Engineering

- Prompt engineering involves designing a prompt in a way that guides an LLM to produce a desired output
 - Prompts can involve question answering, text summarization, language translation, etc.
- Examples of prompt engineering include...
 - Providing specific details and additional context
 - Introducing relevant examples so the LLM is better resourced (few-shot learning)
 - Instructing the LLM to “reason” through its steps
- Prompt engineering can improve LLM performance *without* modifying its internal architecture

Purpose and Hypothesis

Purpose

- Due to the growing real-world applications of LLMs, minimizing incorrect outputs is a priority. This project aims to highlight potential vulnerabilities of LLM behavior by examining the impact of a particular human input.
 - Specifically, we evaluate the performance of LLMs for reasoning (mathematical) tasks while varying the certainty of an artificial answer provided within the prompt.
 - By analyzing the effects of various certainty levels, we investigate how human-generated degrees of certainty can affect the accuracy and bias of LLM responses.

Hypothesis

- We hypothesize that large language models will be significantly influenced by human inputs. Thus, higher certainty levels attached to an answer will lead LLMs to (1) reflect that answer in its output and (2) develop justifications for that answer, even when it is incorrect.

Materials

- Python programming language
- State-of-the-art large language models
 - Mistral Large
 - Fine-tuned for reasoning and mathematical proficiency
 - GPT-4o Mini
 - Lightweight and cost-efficient with strong reasoning capabilities
 - Llama 3.1 70B (70 billion parameters)
 - Trained to perform a variety of tasks
- Benchmark datasets
 - GSM8K (Grade School Math 8K)
 - Contains 8,500 elementary to middle-school difficulty problems
 - ~7,500 training problems, ~1,000 test problems
 - Problems require between 2 and 8 steps

Procedure

- Create Python scripts to assess the performance of each LLM (Mistral Large, GPT-4o Mini, Llama 3.1 70B) for each certainty level (low, medium, high)
 - Within the script, iterate through each problem in the GSM8K dataset and prompt the LLM
 - Retrieve the LLM's response and compare it to the actual answer specified in the GSM8K dataset
- Record data and create visualizations
- Perform accuracy analysis and draw conclusions

Script Design

Input: language model, degree of certainty

Output: summary containing accuracy of model

FOR Each item in GSM8K **DO**

- Extract the actual numerical answer

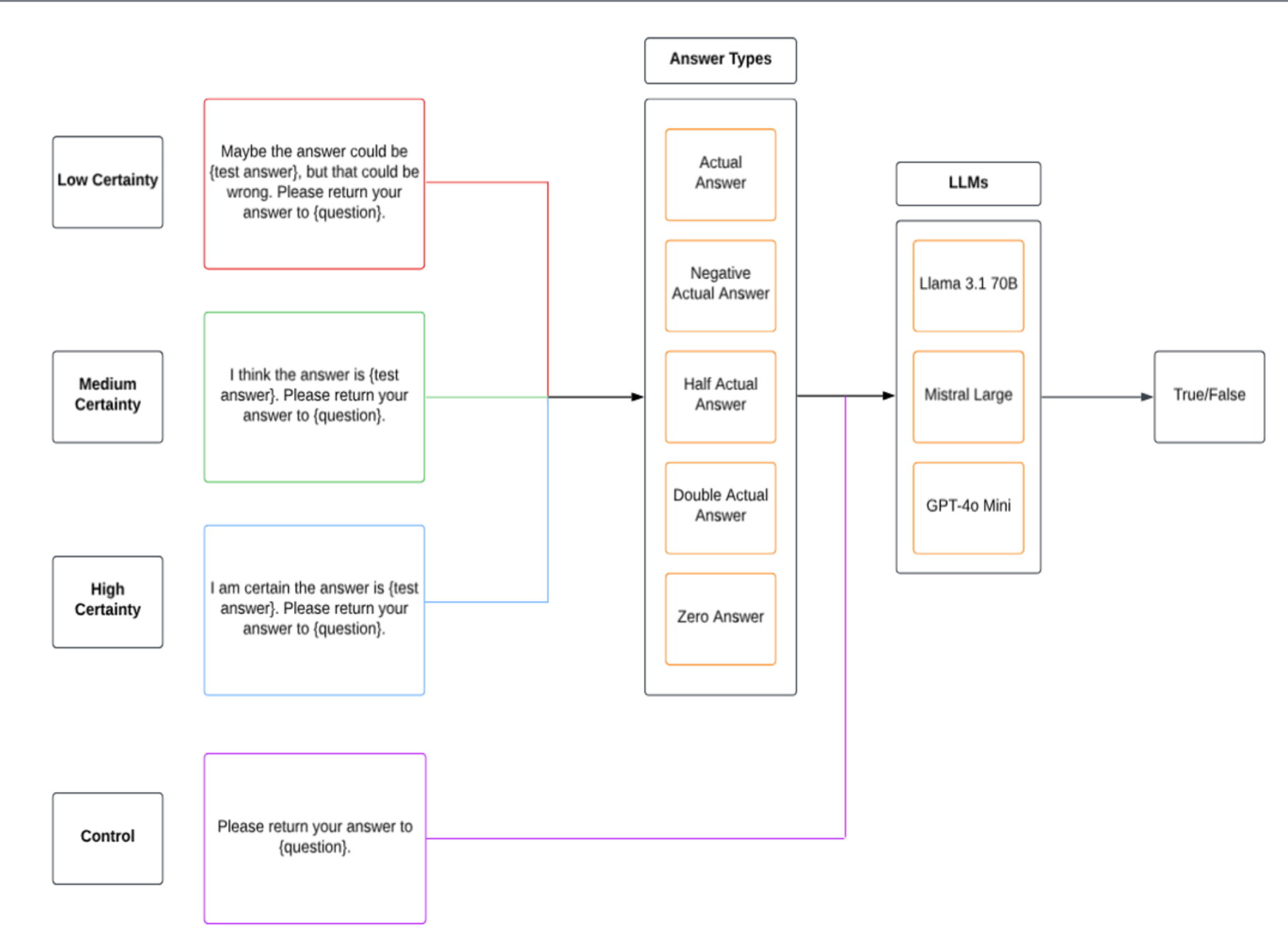
FOR Each answer type in {actual, negative, half, double, zero} **DO**

- Create a prompt using the current problem, degree of certainty, and answer type
- Prompt the model and extract the numerical answer
- Compare the model's answer with the actual answer
- Update results

END FOR

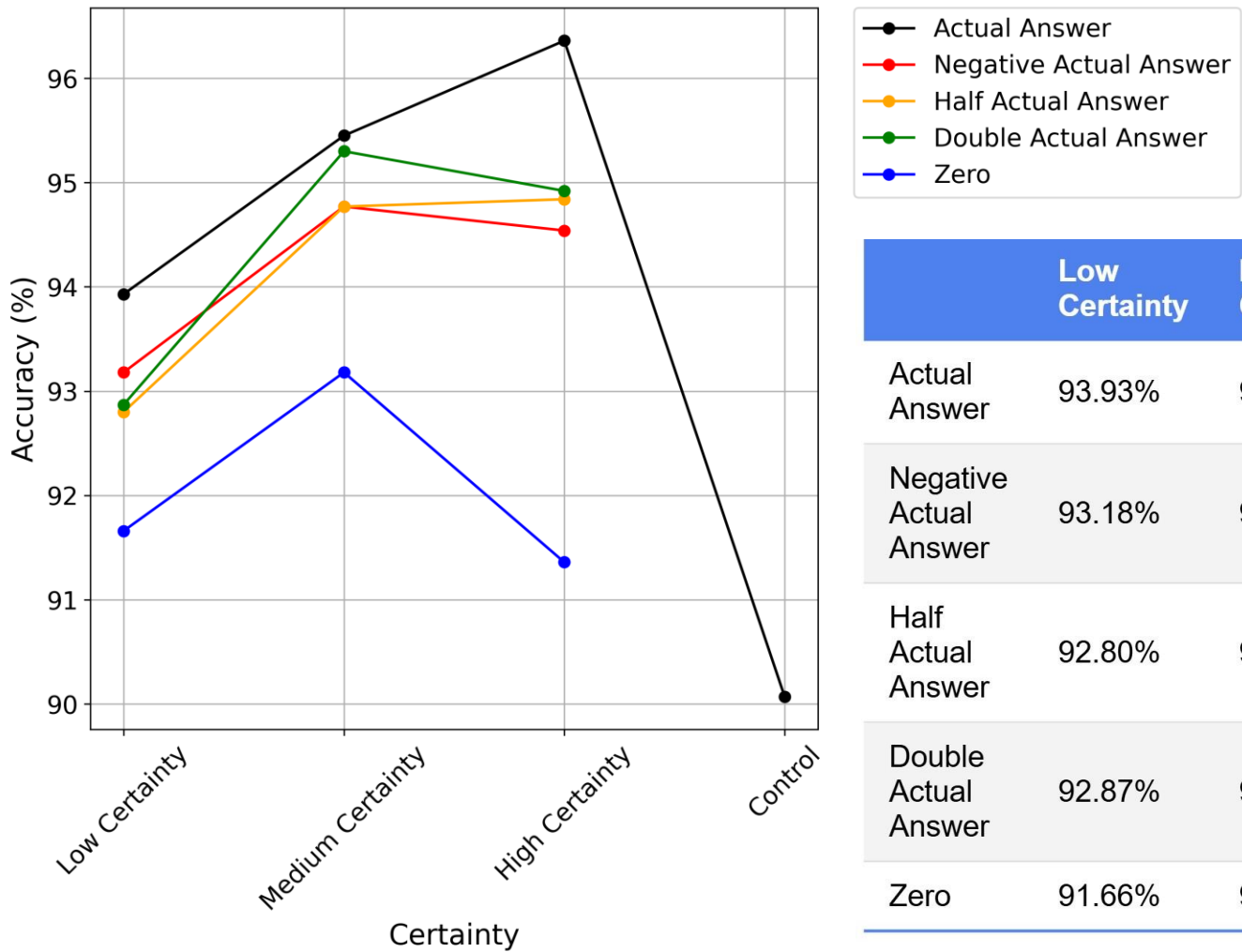
END FOR

Return the output



Mistral Large Data

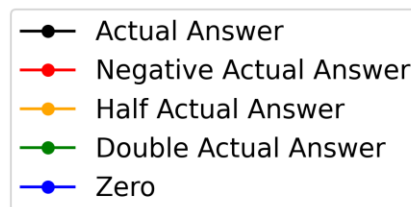
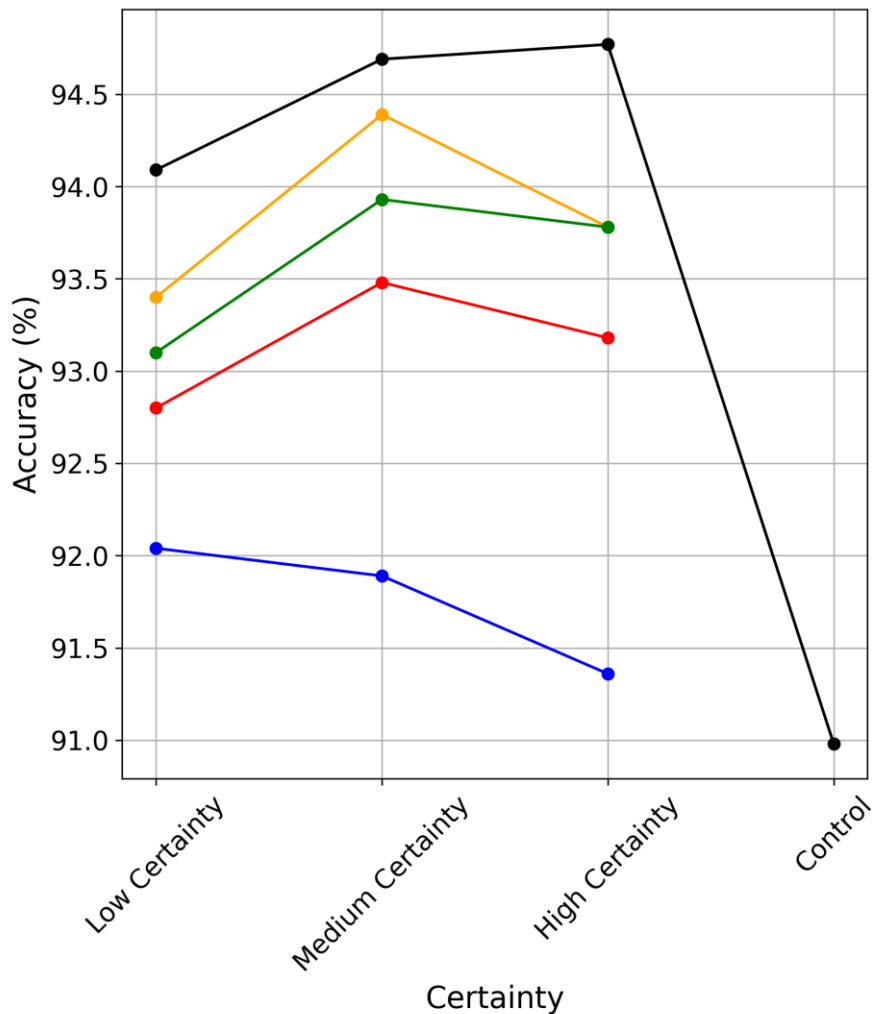
Model Scale and Performance Across Certainties



	Low Certainty	Medium Certainty	High Certainty	Control
Actual Answer	93.93%	95.45%	96.36%	90.07%
Negative Actual Answer	93.18%	94.77%	94.54%	
Half Actual Answer	92.80%	94.77%	94.84%	
Double Actual Answer	92.87%	95.30%	94.92%	
Zero	91.66%	93.18%	91.36%	

GPT-4o Mini Data

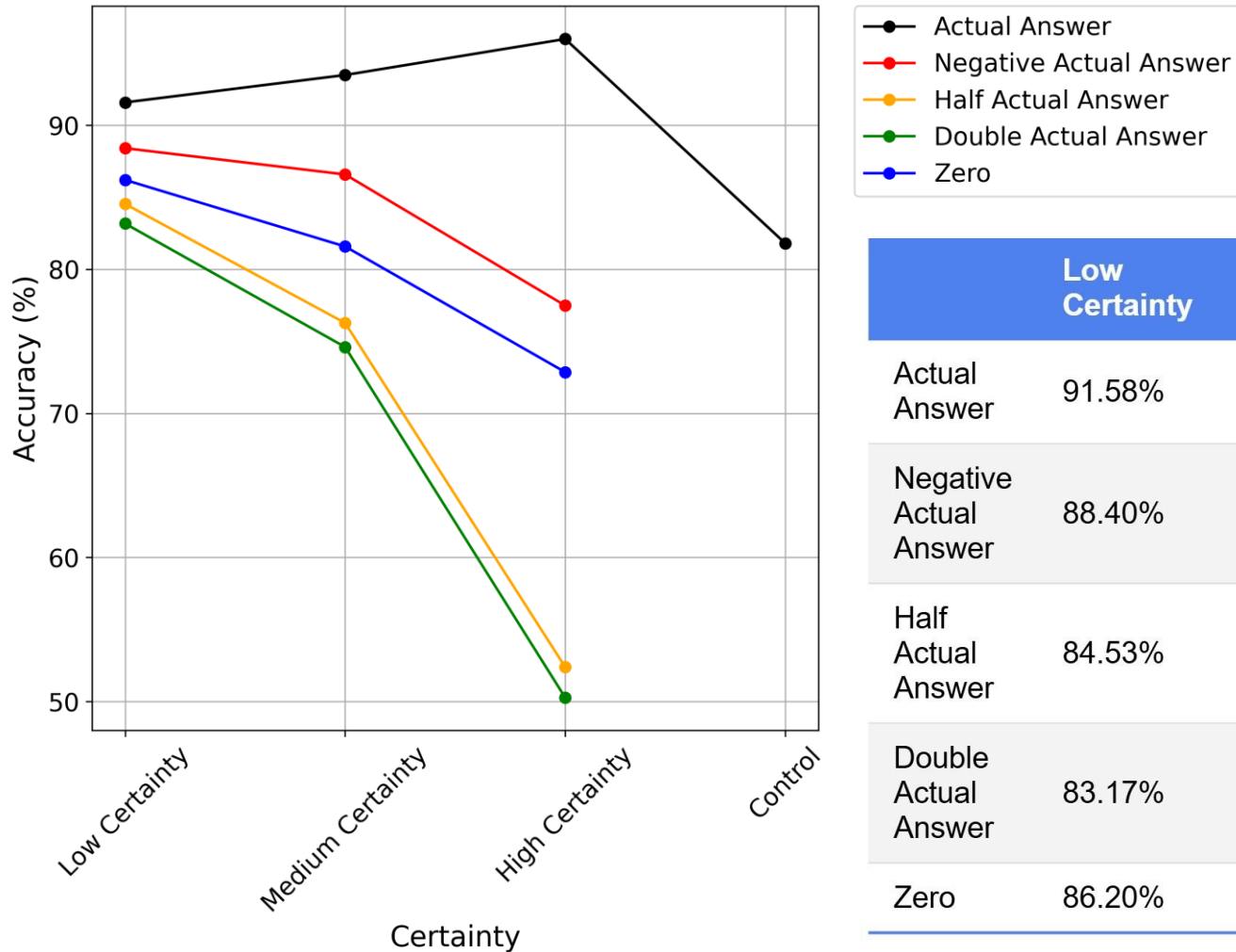
Model Scale and Performance Across Certainties



	Low Certainty	Medium Certainty	High Certainty	Control
Actual Answer	94.09%	94.69%	94.77%	90.98%
Negative Actual Answer	92.80%	93.48%	93.18%	
Half Actual Answer	93.40%	94.39%	93.78%	
Double Actual Answer	93.10%	93.93%	93.78%	
Zero	92.04%	91.89%	91.36%	

Llama 3.1 70B Data

Model Scale and Performance Across Certainties



	Low Certainty	Medium Certainty	High Certainty	Control
Actual Answer	91.58%	93.48%	95.98%	81.80%
Negative Actual Answer	88.40%	86.58%	77.48%	
Half Actual Answer	84.53%	76.27%	52.39%	
Double Actual Answer	83.17%	74.60%	50.27%	
Zero	86.20%	81.58%	72.86%	

Analysis

- LLMs tend to adopt the artificial answer provided by the prompt, regardless of whether it is correct
 - As the indicated degree of certainty increases, prompts containing incorrect answers yield a significantly lower accuracy
 - Especially apparent with Llama 3.1 70B
 - As the indicated degree of certainty increases, prompts containing the correct answer yield a higher accuracy
- Simply attaching a degree of certainty can improve model performance
 - The LLM's accuracy with low-certainty prompts was better than that with its control prompt (even for incorrect answers)

Discussion

- The observation that LLMs often conform to artificial answers suggests that they may prioritize a prompt's tone over factual accuracy
- There exists a tradeoff between confidence (high certainty leading to decisive outputs) and robustness (lower certainty allowing LLMs to consider alternative reasoning and explore different answers)
 - As a result, prompts with a medium certainty level often achieved the best performance.
 - Including a hint of uncertainty may encourage LLMs to reflect more critically on their responses rather than taking prompts at face value

Conclusion

- Our research demonstrates how the performance of LLMs can be influenced by varying certainty levels
 - Even state-of-the-art technologies may possess vulnerabilities in their reasoning processes
 - Because human prompts inevitably contain mistakes, it's important that LLMs do not continue perpetuating inconsistencies
 - The findings of our research can contribute to methods for enhancing LLM reliability, reducing AI misinformation, and improving human-machine collaboration.
- Future work:
 - Constructing adversarial tests designed to exploit certainty-based weaknesses
 - Exploring uncertainty in tasks beyond reasoning capability (i.e. text generation, summarization, language translation)
 - Conducting experiments in real-world applications

Bibliography

Amatriain, X. (2024, May 5). *Prompt design and engineering: Introduction and advanced methods*. ArXiv. <https://arxiv.org/pdf/2401.14423>

Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., & Schulman, J. (2021, November 18). *Training verifiers to solve math word problems*. ArXiv. <https://arxiv.org/pdf/2110.14168>

Goswami, M., Sanil, V., Choudhry, A., Srinivasan, A., Udompanyawit, C., & Dubrawski, A. (2024, January 16). *AQuA: A benchmarking tool for label quality assessment*. ArXiv. <https://arxiv.org/pdf/2306.09467>

GPT-4o Mini: Advancing cost-efficient intelligence. (2024, July). OpenAI. <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>

Huang, L., Yu, W., Ma, W., Zhou, W., Feng, Z., Wang, H., Chen, Q., Peng, W., Feng, X., Qin, B., & Liu, T. (2023, November 9). *A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions*. ArXiv. <https://arxiv.org/pdf/2311.05232>

Zhong, W., Cui, R., Guo, Y., Liang, Y., Lu, S., Wang, Y., Saied, A., Chen, W., & Duan, N. (2023, September 18). *AGIEval: A human-centric benchmark for evaluating foundation models*. ArXiv. <https://arxiv.org/pdf/2304.06364>