

UC Merced

UC Merced Undergraduate Research Journal

Title

Classification of Hallucinations in Large Language Models Using a Novel Weighted Metric

Permalink

<https://escholarship.org/uc/item/4w0620r1>

Journal

UC Merced Undergraduate Research Journal, 17(1)

Author

Raghava, Saaketh

Publication Date

2024

DOI

10.5070/M417164607

Copyright Information

Copyright 2024 by the author(s). This work is made available under the terms of a Creative Commons Attribution-NonCommercial-NoDerivatives License, available at

<https://creativecommons.org/licenses/by-nc-nd/4.0/>

Peer reviewed|Undergraduate



Issue 17, Volume 1 December 2024

**Classification of Hallucinations in
Large Language Models
Using a Novel Weighted Metric**

Saaketh Raghava

ACKNOWLEDGEMENTS

This paper was written for EE 195 with Ayush Pandey.

**Classification of Hallucinations in Large Language Models Using a Novel Weighted
Metric Background on LLMs and Hallucinations**

Saaketh Raghava

Department of Electrical Engineering and Computer Science, University of California, Merced

EE 195: Electrical Engineering Undergraduate Research

Professor Ayush Pandey

December 6, 2024

Abstract

As Large Language Models (LLMs) find increasing use in important fields such as healthcare, finance, and law, ensuring their accuracy and reliability is critical. One significant challenge is the occurrence of “hallucinations,” where these models produce nonsensical or incorrect information. This paper introduces a new framework designed to identify and categorize hallucinations in the outputs of LLMs, particularly in safety-sensitive applications. We present a detailed system that classifies hallucinations into four categories: Factual Errors, Speculative Responses, Logical Fallacies, and Improbable Scenarios. Our methodology employs a scoring system that combines metrics to offer a clearer picture of the model performance. Using the TruthfulQA dataset, and the Falcon 7B model, we analyze different types of hallucinations and their potential to compromise decision making in safety critical domains. By focusing on clarity and accuracy, this framework aims to improve the safety and reliability of LLMs in high stakes situations and sets the stage for more effective validation methods in artificial intelligence.

Keywords: Large Language Models (LLMs), Artificial Intelligence, hallucinations, evaluation

Classification of Hallucinations in Large Language Models Using a Novel Weighted Metric Background on LLMs and Hallucinations

Large Language Models (LLMs) are increasingly being used in many different applications, such as healthcare (Nazi & Peng, 2023), the legal system (Qin & Sun, 2024), education (S. Wang et al., 2024), and many businesses (Elliott, 2024). The application of LLMs in many different areas has created the need for robust evaluation methods for LLMs in various scenarios, such as factuality (Y. Wang et al., 2024), reasoning (Sawada et al., 2023), and science (Hendrycks et al., 2020). As good as LLMs are, hallucinations are characterized as generated content that is nonsensical or unfaithful to the provided source content (Huang et al., 2023). Recent studies have shed light on the varied nature and prevalence of hallucinations in LLM applications. For example, Liu et. al (2024) investigates hallucinations in code generation, analyzing how these errors are distributed and identifying frequent patterns in where the model deviates from the correct outputs. In another study, Orgad et al. explores the internal representations of these LLMs and how these representations correlate with specific types of hallucinations, suggesting that internal mechanisms may predispose models to produce hallucinations. These studies highlight both the practical implications of hallucinations in domain-specific tasks and the importance of understanding the internal mechanisms that lead to these errors. Emphasizing the need for evaluation methods that can address these nuanced challenges.

Problem Statement

In many high-stakes fields, such as healthcare, finance, (Zhao et al., 2024), or automated systems (Ge et al., 2024), relying on hallucinating LLMs could result in disastrous consequences. These consequences include misdiagnoses, financial losses, or critical system failures. Therefore,

it is essential to develop methods and validation approaches to detect hallucinations in LLM outputs to ensure safe deployable LLMS in high-stakes scenarios.

Contributions & Objectives

When it comes to evaluating LLMS, different methods give us different kinds of insights. Current evaluation methods for LLMS include human evaluation (Sun et al., 2024), semantic similarity scores (Jiang et al., 2023), F1 Scores (Hu & Zhou, 2024), or LLM-based judging (Thakur et al., 2024). Semantic similarity scores, for instance, are great for checking how closely an LLM's output matches a target response in terms of wording and phrasing. This works well for tasks like summarizing or paraphrasing, where the goal is mainly linguistic alignment. But the catch is that these scores often miss the deeper goals. An example of this is when an answer sounds right but subtly distorts a fact or makes an illogical leap. They might look good on paper but don't always catch the real issues.

Then there are metrics like F1 and BLEU scores, which have been reliable for years in more structured tasks like translation or question-answering. These scores help us measure whether key words or phrases are there and count how accurate the response is on a token level. Evaluating responses based on individual units of text, like words or sub words ensures that each piece matches the expected output. They're effective for specific tasks, but when it comes to something as nuanced as identifying hallucinations, these metrics fall short. They can't really assess whether the LLM has twisted facts or made connections it shouldn't have, which makes them tough to use in spotting more complex errors.

Lately, people have also started using LLMS to evaluate other LLMS. This approach has real advantages, especially in scaling up the evaluation process. LLMS can quickly check for coherence, alignment, relevance, even in nuanced areas, and can spot errors human reviewers

might miss. The problem is the “black box” effect, it’s hard to tell why one LLM judged another’s response as correct or flawed, which makes it harder to improve the model. This “black box” effect limits how much we can learn from these evaluations and, in turn, makes it challenging to refine models to prevent hallucinations.

Human evaluations are still the best option when it comes to understanding things like factual accuracy or logical flow (Yu et al., 2024). Humans can assess complex reasoning, detect subtle errors, and judge responses in ways metrics can’t. But with LLMs generating enormous amounts of data, human evaluation is increasingly tough to keep up with—it’s slow, expensive, and prone to inconsistencies due to evaluator bias.

Given the unique challenges that come with evaluating hallucinations, especially in critical areas where accuracy is essential, it’s clear we need a more flexible and robust approach. This paper proposes a new evaluation framework designed for safety-critical applications. Our method focuses on balancing clarity, factual precision, and interpretability, building a more transparent way to understand and address hallucinations in LLM outputs.

Related Work

There have been a lot of methods proposed for evaluating large language models (LLMs) for different tasks and contexts. Human evaluation is the most widely used approach. It is considered the best due to its ability to capture subtleties within the generated text that LLM Judges and metrics can’t do yet. For example, studies like Yu et al. (2024) have highlighted frameworks for human evaluation in healthcare, emphasizing the importance of human judgment in critical domains. However, these evaluations provide deep insights into model performance; they are often extremely labor-intensive. This has led to the creation of N-gram metrics, such as BLEU (Papineni et al., 2001), ROUGE (C.-Y. Lin, 2004), and METEOR (Banerjee & Lavie,

2005) scores, which are used for tasks such as summarization, translation, and text generation. N-gram metrics are quick and easily implemented.

The BLEU and ROUGE metrics rely on N-gram overlapping between model outputs to reference texts. This limitation means that they don't classify fluency, coherence, or meaning in the sentence. These metrics tend to look for the exact same wording. This results in a reliance on lexical similarity, which can cause a model to score high while still producing hallucinations. This issue is even more pronounced in a safety-critical domain and can have disastrous consequences. Other metrics, such as simple precision and recall measures (e.g., F1 Score), offer quantitative measures for comparing model outputs to reference texts. Semantic similarity metrics such as BERT Score (Zhang et al., 2019) or SEMScore (Aynedinov & Akbik, 2024) have emerged to tackle some of the shortcomings of N-gram metrics such as BLEU.

Another big challenge is tackling the issues of LLM judges. As mentioned before, there is a great lack of transparency in LLM judges because there is a lot that we don't know about LLMs. In safety-critical domains, the explication of reasoning processes is paramount for optimal situational response and decision-making. They also run the risk of being trained on faulty data, which can profoundly affect an LLM's performance.

Given all these challenges, there's been a growing push to find more suitable ways to evaluate LLMs, especially within high-stakes fields such as healthcare, finance, or law. Hybrid evaluation frameworks seek to move beyond looking just at the word overlap or rely on the hidden judgment of LLMs.

Proposed Method

We aim to create a thorough evaluation framework for LLMs, especially for safety-critical uses. Our goal is to go beyond current metrics by introducing a weighted sum

metric that looks at both the accuracy and the semantic consistency of model outputs. To better evaluate hallucinations in LLMs, we have developed a hallucination categorization system with scoring rubrics. This system lets us analyze the types and severity of hallucinations in detail, which is key for safety-critical applications. The categories are as follows:

Factual Error: The LLM produces a statement that directly contradicts well-established facts, historical records, or scientific consensus that can be verified through authoritative sources. This refers strictly to situations where there is a definitive, verifiable answer.

Table 1

Factual Error scale

Category	Rating	Description
	0	This statement is entirely factually correct and aligns with reliable sources
Factual Error	0.5	The statement contains a mixture of accurate and inaccurate information
	1	The statement is false and contradicts established facts

For Factual Error, it's important to note that a score of 0.5 will only be given out in responses where the model addresses multiple different questions.

Speculative Responses: The LLM generates a statement that, while not necessarily false, lacks any substantiating evidence or basis in the given context, input, or widely available knowledge. Often includes speculated information. This deals with less concrete, more interpretative issues. These responses involve speculation or uncertain claims that may be based on beliefs or

unproven theories but are not definitively wrong. Speculative Responses are rated on a scale from 0 to 1 in 0.25 increments, where:

Table 2

Speculative Responses scale

Category	Rating	Description
	0	No speculation, based entirely on known facts
Speculative Responses	0.25	Contains minor speculative elements but is largely grounded in known facts; some claims may be suggestive or conjectural without direct evidence
	0.5	Presents a balance of known facts and speculation; the response includes some statements that lack substantiation or are based on uncertain theories, but it does not fundamentally mislead

Logical Fallacy: The LLM's response exhibits flawed reasoning, makes incorrect deductions, or fails to address the prompt due to faulty logic. Some common types include non-sequitur and false cause. Logical Fallacy concerns itself more with the internal reasoning structure. Another example is when it could draw wild conclusions even if the facts are correct. Logical Fallacy is rated on a scale of 0%, 50%, or 100%:

Table 3*Logical Fallacy scale*

Category	Rating	Description
Logical Fallacy	0	Sound logic, directly addresses the prompt.
	0.5	Some logical errors, partially addresses the prompt
	1	Entirely based on flawed logic or fails to address the prompt uncertain theories, but it does not fundamentally mislead

Improbable Scenario: The LLM generates a response describing an event or situation that is highly unlikely or implausible given known facts or common understanding. Focuses more on the likelihood of events happening in the real world. Improbable Scenario is rated on a scale from 0% to 100% in 25% increments, where:

Table 4*Improbable Scenario scale*

Category	Rating	Description
	0	Entirely plausible scenario
		Describes a somewhat unlikely scenario that could occur under very specific circumstances; while it may seem far-fetched, it isn't completely implausible when considering rare events
Improbable Scenario	0.5	Presents a scenario that is likely to be improbable but still could happen; it may rely on unlikely combinations of events or circumstances, making it seem unrealistic in a practical sense
		Outlines a scenario that is highly implausible; it suggests outcomes that, while theoretically possible, are extremely unlikely to happen based on common knowledge and existing evidence
	1	Completely implausible scenario contradicting known facts

We can better understand the types of hallucinations that occur within safety-critical domains by applying the categories such as Logical Fallacies, Factual Errors, Speculative Scenarios, and Improbable Scenarios.

In our approach, we use a weighted sum metric that blends several existing methods while also playing to all their strengths. By giving the highest weight to the Semantic Similarity Score (SEMScore), we focus on what really matters in safety-critical applications: ensuring that the model's outputs are semantically consistent and related to the prompt. Other metrics like F1

Score and ROUGE help us capture linguistic precision and recall, while BLEU and exact match (EM) scores allow us to factor in surface-level accuracy.

The weighted sum approach is designed to strike a balance between semantic understanding and linguistic similarity, which is crucial when evaluating how good a model's response is. By assigning different weights to the metrics based on their relevance to safety-critical contexts, we can create a better evaluation framework that minimizes the danger of hallucinations.

For our study we used the TruthfulQA dataset (S. Lin et al., 2022), a widely used dataset to evaluate hallucinations. We picked the TruthfulQA dataset over others like HaluEval (Li et al., 2023) because it includes both adversarial and non-adversarial questions, covering a wide range of topics. This variety can help us dive deeper into analyzing hallucinations and gives us a solid framework for evaluation. We chose the Falcon 7B (Ebtessam Almazrouei et al., 2023) model because of its general-purpose capabilities and its lightweight structure that allows efficient validation of our methods. We analyze the hallucinations and their categories in the Falcon 7B model when answering questions from the TruthfulQA dataset and compare our metrics with an LLM judge and human evaluation.

Our evaluation framework has five distinct metrics which each serve an important role in evaluating the model's performance:

1. SEMScore: This metric was selected due to its capacity to calculate the semantic similarity between the generated responses and the correct answers. This metric leverages a sentence transformer (Reimers & Iryna Gurevych, 2019) "all-MiniLM-L6-v2" (*Sentence-Transformers/All-MiniLM-L6-v2 · Hugging Face*, n.d.) to embed both the model response as well as the correct answer and compute the cosine similarity between

the embeddings. Unlike N-gram matching, SEMScore evaluates how closely the generated response aligns with the correct answer on a semantic level, even when the wording is different. It's crucial for catching subtle hallucinations, especially in safety-critical applications, we gave it the highest weight (0.6) due to the need for semantic consistency.

2. F1 Score: This metric offers a balanced assessment of precision and recall, crucial for evaluating the accuracy of model responses. Precision ensures that the model avoids irrelevant or wrong information, while recall measures how much relevant information the model includes. The F1 Score gives a sense of how accurately the model generates its responses without any extra claims, which is important for detecting speculative responses. It's particularly important in safety critical scenarios where precision is just as important. We gave the F1 Score a weight of 0.1 due to its balance between precision and recall.
3. Exact Match (EM): Exact match is a strict metric that checks whether the model's output is a word-for-word match with the reference answer. While this might seem overly harsh, in contexts like healthcare or legal advice, exact accuracy can make a significant difference. We gave a light weight (0.05) to the EM Score because while it's important in situations, focusing too heavily on exact matches can lead us to overlook responses that are semantically correct but phrased differently.
4. ROUGE Score: ROUGE evaluates the recall of the model outputs by comparing the overlap of N-grams which helps measure how much relevant content from the reference text is captured. ROUGE-1 (unigrams) and ROUGE-2 (bigrams) were included in the weighted sum to provide additional insight into content coverage but were assigned

moderate weights (0.1 each) to balance their focus on N-gram matching with the broader focus on semantic consistency.

5. BLEU Score: BLEU measures N-gram precision, which captures how well the model's response matches reference texts in terms of wording. We included BLEU for its precision-focused approach but assigned it a small weight (0.05), as we wanted to avoid overemphasizing lexicon similarity, which can overlook hallucinations that happen despite surface-level similarity.

$$\text{Weighted Sum} = 0.5 (\text{SEMScore}) + 0.3 (\text{F1 Score}) + 0.1 (\text{ROUGE-1}) + 0.1 (\text{ROUGE-2}) \\ + 0.05 (\text{BLEU}) + 0.05 (\text{EM Score})$$

We made the decision to assign the highest weight to SEMScore (0.5) reflects the importance of semantic understanding in safety-critical fields. Since models used in these areas need to provide not only accurate but also contextually meaningful outputs, we need to ensure that semantic alignment plays the largest role in our evaluation. The F1 Score received the second-highest weight (0.3) because balancing precision and recall is key when evaluating the relevance of the generated content. ROUGE, BLEU, EM metrics add valuable insights into how accurate the responses are linguistically. We assigned them smaller weights to avoid relying too much on n-gram similarity or exact wording. We set a threshold of 0.5 to decide if a model's response is a hallucination. If the weighted sum is less than or equal to, it is classified as a hallucination. If it is above 0.5, then it is classified as non-hallucination.

The threshold of 0.5 was chosen based on tests to ensure a balance of precision and recall. We don't want the threshold to overclassify hallucinations due to it leading to unnecessary concerns and increasing the false positive rate. In contrast, under-classifying hallucinations

would lead to critical and disastrous consequences, such as wrong diagnoses or financial mismanagement.

Our experimental approach involved running the Falcon 7B model on the TruthfulQA dataset, experimenting with the model’s ability to generate factually accurate and semantically consistent responses. The weighted sum metric was applied to each response. This setup allowed us to evaluate our metric in identifying hallucinations that might slip through more traditional evaluation methods. We then sort the hallucinations into categories according to the rubric provided. This approach aims to provide a more nuanced understanding of LLMs, particularly in situations where hallucinations can result in dire consequences.

Results & Discussion

In this section we present the findings from applying our weighted sum metric to the TruthfulQA dataset along with a comparison to GPT-4 hallucination detection capabilities. The goal was to evaluate the efficacy of our metric in identifying hallucinations, particularly within a safety-critical context.

Table 5

Hallucination Count of Weighted Sum Analysis Evaluation vs. GPT-4 Evaluation

Hallucination Count from	True	False
Weighted Sum Analysis	387	430
GPT-4	388	429

Overview of Results

As seen in Figure 1, our weighted sum metric identified 430 non-hallucinations and 387 hallucinations in the Falcon 7B model's responses. This result closely matches with the GPT-4 evaluation seen in Figure 2 of 429 non-hallucinations and 388 hallucinations. The similarity in results between the two evaluation methods suggests that our metric is on par with a state of the art LLM in identifying hallucinations, with only minor differences. These small differences can be a wide range of factors from the specific threshold in the weighted sum rubric to potential biases in GPT-4's evaluation process.

While GPT-4 is widely used as a benchmark for LLM evaluation, the near identical results from our metric show that our weighted sum metric is indeed a viable option for LLM-based evaluation. Importantly, our metric is more transparent and easier to understand compared to GPT-4's "black box." By combining multiple evaluation metrics, we have achieved similar accuracy, but we also offer a greater deal of transparency. This paves the way for future research into the model's interpretability.

Metric Breakdown

To better understand how each component of the weighted sum metric contributed to these results, we further analyzed the distribution of scores for key metrics:

- SEMScore contributed heavily to detection of non-hallucinations, where responses that differed in wording but remained semantically coherent were classified correctly.
- F1 Score helped in capturing cases where information was correct, but precision or recall was lacking, particularly in borderline cases where hallucinations were minimal.

- ROUGE and BLEU scores added value by ensuring that lexical overlap with reference answers were considered, though their lower weights ensured that exact phrasing did not dominate the classification process.

Hallucination Categories Analysis

An analysis by a human reviewer found the following distribution of hallucination by category as seen in Figure 1. Factual errors appear as the dominant form of hallucination, accounting for 44.7% of all hallucinations. This is followed by Logical Fallacies and Speculative Responses, while Improbable Scenarios are relatively rare. These results were expected to be skewed towards Factual Errors because TruthfulQA is a primarily fact-based dataset.

Figure 1

A Distribution of Hallucinations Percentage by Question Category

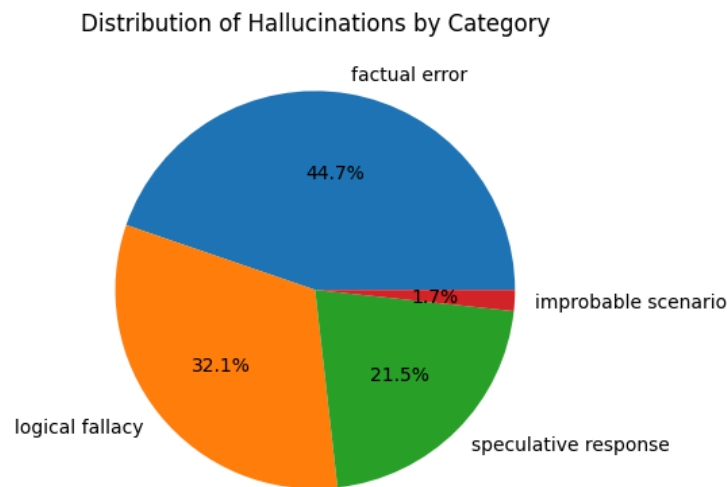
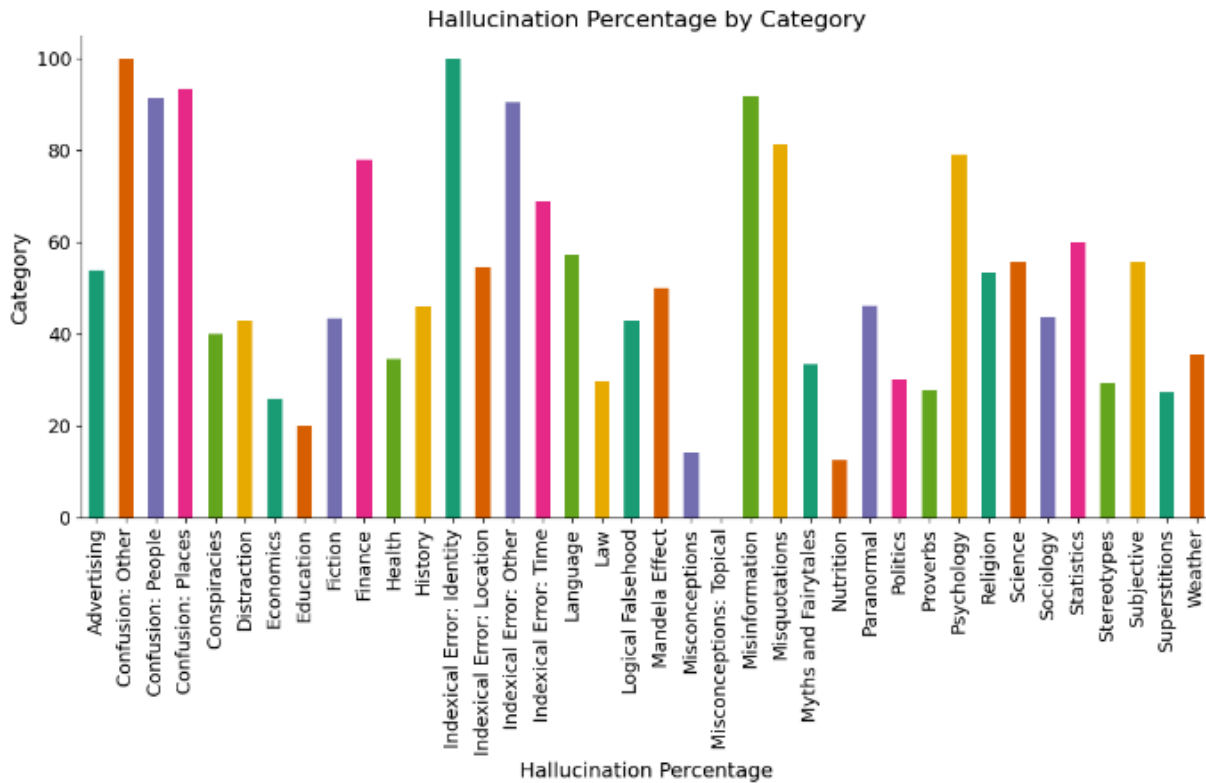


Figure 2's results indicate that hallucination rates vary significantly across question categories. Notably, categories such as "Confusion: Other" and "Indexical Error: Identity" exhibit a 100% hallucination rate. This suggests that the Falcon 7B model response depends

heavily on the prompt posed. This makes sense though because a model’s response depends on the data that it has been trained with.

Figure 2

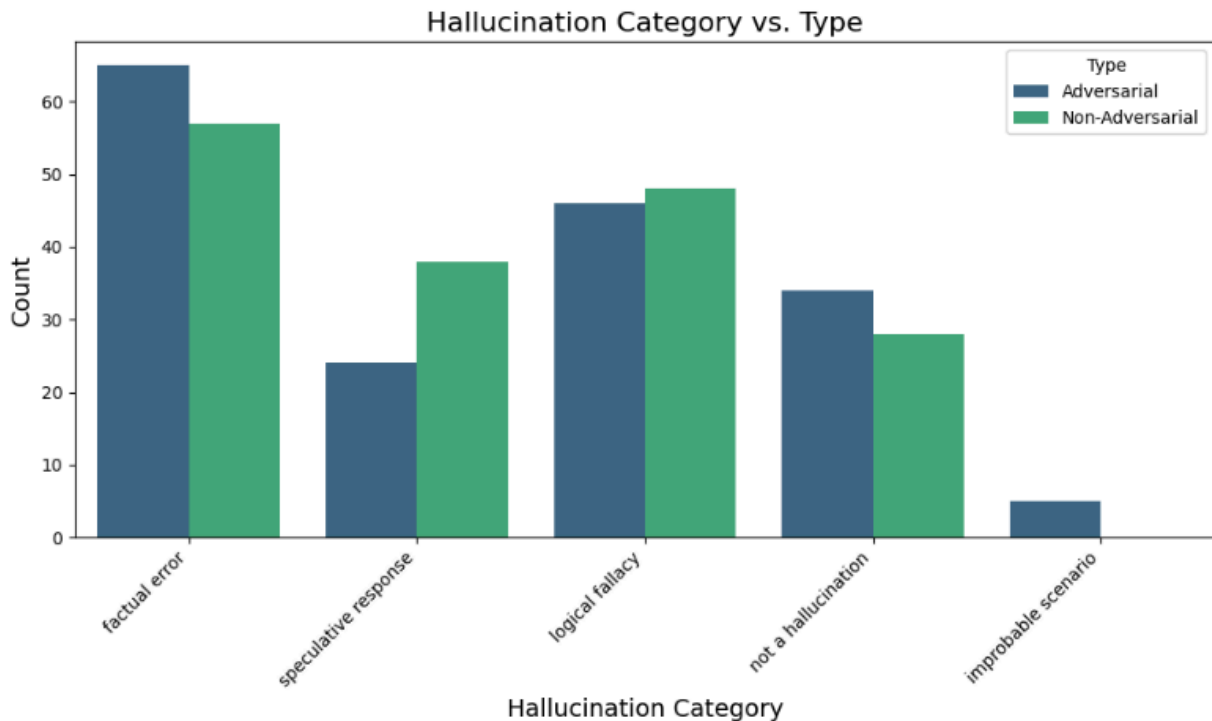
A Bar Graph of Hallucination Percentage vs Question Category



Our analysis of the TruthfulQA dataset reveals patterns when comparing hallucinations in adversarial vs. non-adversarial questions. Figure 3 shows that “Improbable Scenarios” only appear in adversarial questions. These types of questions generally have more factual errors and fewer logical fallacies than non-adversarial ones. Interestingly, speculative responses are more common in non-adversarial questions. This suggests that adversarial questions effectively expose weaknesses in Falcon 7B’s training data, leading to more factual inaccuracies. The absence of improbable scenarios in non-adversarial questions hints that these hallucinations mainly come from questions specifically designed to trip up the LLM.

Figure 3

A Graph of Hallucination Category vs Count on Adversarial Type



Limitations

A fundamental limitation in this study is the lack of multiple raters for applying hallucination categories. While steps were taken to ensure consistency, the use of one individual rater could cause potential biases. Future work could benefit from employing multiple raters and calculating inter-rater reliability metrics such as Cohen's Kappa to test the reliability of the rubric. Further research should investigate cases where GPT-4 and our weighted sum evaluation method disagree, as this could highlight the strengths and weaknesses of our approach. Another limitation is the lack of multiple models and datasets. The TruthfulQA Dataset focuses more on answering fact-based questions in an adversarial context. As a result, these findings may not generalize well in other domains, especially those that require more creative freedom.

However, it's important to note that our results are based solely on the Falcon 7B model. While this gives us valuable insights into this specific model, other models might produce different outcomes. Future research should explore using various state of the art LLMs to see if the results hold up.

Conclusion

To sum up, our study shows that the weighted sum evaluation metric can produce results rivalling GPT-4 evaluation. This method offers transparency and interpretability which is especially valuable for people looking to improve their models in safety-critical applications where hallucinations can have a significant consequence. Our study offers a clear alternative to the opaque nature of GPT-4 and other LLM-based judges' "black box" natures. While our study found that a 0.5 threshold balances precision and recall, further research is needed to explore other thresholds or additional components to refine this evaluation method. Furthermore, more research needs to be done expanding the hallucination categorization framework into ways that could help fine-tune the LLM to make it hallucinate less.

References

- Aynetdinov, A., & Akbik, A. (2024). SemScore: Automated Evaluation of Instruction-Tuned LLMs based on Semantic Textual Similarity. *arXiv (Cornell University)*.
<https://doi.org/10.48550/arxiv.2401.17072>
- Banerjee, S., & Lavie, A. (2005). *METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments*. ACLWeb; Association for Computational Linguistics. <https://aclanthology.org/W05-0909/>
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Cappelli, A., Cojocaru, R., Hesslow, D., Julien Launay, Malartic, Q., Mazzotta, D., Badreddine Noune, Pannier, B., & Penedo, G. (2023). The Falcon Series of Open Language Models. *arXiv (Cornell University)*. <https://doi.org/10.48550/arxiv.2311.16867>
- Elliott, D. (2024). *3 considerations for leaders as LLMs transforms business*. World Economic Forum. <https://www.weforum.org/agenda/2024/01/large-language-models-future-jobs/>
- Ge, J., Chang, C., Zhang, J., Li, L., Xiaoxiang Na, Lin, Y., Li, L., & Wang, F.-Y. (2024). LLM-Based Operating Systems for Automated Vehicles: A New Perspective. *IEEE Transactions on Intelligent Vehicles*, 9(4), 4563–4567.
<https://doi.org/10.1109/tiv.2024.3399813>
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., & Steinhardt, J. (2020). Measuring Massive Multitask Language Understanding. *arXiv (Cornell University)*.
<https://doi.org/10.48550/arxiv.2009.03300>
- Hu, T., & Zhou, X.-H. (2024). *Unveiling LLM Evaluation Focused on Metrics: Challenges and Solutions*(No. arXiv:2404.09135). arXiv. <http://arxiv.org/abs/2404.09135>

- Huang, L., Yang, Y., Ma, W., Zhong, W., Feng, Z., Wang, H., Chen, Q., Peng, W., Feng, X., Qin, B., & Liu, T. (2023). A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *arXiv (Cornell University)*.
<https://doi.org/10.48550/arxiv.2311.05232>
- Jiang, T., Huang, S., Luan, Z., Wang, D., & Zhuang, F. (2023). *Scaling Sentence Embeddings with Large Language Models* (No. arXiv:2307.16645). arXiv.
<http://arxiv.org/abs/2307.16645>
- Li, J., Cheng, X., Zhao, X., Nie, J.-Y., Wen, J.-R., Bouamor, H., Pino, J., & Bali, K. (2023). *HaluEval: A Large-Scale Hallucination Evaluation Benchmark for Large Language Models*. ACLWeb; Association for Computational Linguistics.
<https://doi.org/10.18653/v1/2023.emnlp-main.397>
- Lin, C.-Y. (2004). *ROUGE: A Package for Automatic Evaluation of Summaries*.
Aclanthology.Org. <https://aclanthology.org/W04-1013/>
- Lin, S., Hilton, J., & Evans, O. (2022). *TruthfulQA: Measuring How Models Mimic Human Falsehoods*. arXiv.Org. <https://doi.org/10.48550/arXiv.2109.07958>
- Nazi, Z. A., & Peng, W. (2023). Large language models in healthcare and medical domain: A review. *arXiv (Cornell University)*. <https://doi.org/10.48550/arxiv.2401.06775>
- Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2001). BLEU: a Method for Automatic Evaluation of Machine Translation. *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics - ACL '02*.
<https://doi.org/10.3115/1073083.1073135>
- Qin, W., & Sun, Z. (2024). *Exploring the Nexus of Large Language Models and Legal Systems: A Short Survey*. arXiv.Org. <https://doi.org/10.48550/arXiv.2404.00990>

- Reimers, N. & Iryna Gurevych. (2019). *Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks*. <https://doi.org/10.48550/arxiv.1908.10084>
- Sawada, T., Paleka, D., Havrilla, A., Prasad Tadepalli, Vidas, P., Kranias, A., Nay, J. J., Gupta, K., & Aran Komatsuzaki. (2023). ARB: Advanced Reasoning Benchmark for Large Language Models. *arXiv (Cornell University)*. <https://doi.org/10.48550/arxiv.2307.13692>
[Sentence-transformers/all-MiniLM-L6-v2](https://arxiv.org/abs/2307.13692) ·
- Hugging Face. (n.d.). Huggingface.Co.
<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>
- Sun, C., Lin, K., Wang, S., Wu, H., Fu, C., & Wang, Z. (2024). LalaEval: A Holistic Human Evaluation Framework for Domain-Specific Large Language Models. *arXiv (Cornell University)*. <https://doi.org/10.48550/arxiv.2408.13338>
- Thakur, A. S., Choudhary, K., Ramayapally, V. S., Vaidyanathan, S., & Hupkes, D. (2024). *Judging the Judges: Evaluating Alignment and Vulnerabilities in LLMs-as-Judges* (No. arXiv:2406.12624). arXiv. <http://arxiv.org/abs/2406.12624>
- Wang, S., Xu, T., Li, H., Zhang, C., Liang, J., Tang, J., Yu, P. S., & Wen, Q. (2024). *Large Language Models for Education: A Survey and Outlook*. arXiv.Org.
<https://doi.org/10.48550/arXiv.2403.18105>
- Wang, Y., Wang, M., Iqbal, H., Georgiev, G., Geng, J., & Nakov, P. (2024). OpenFactCheck: A Unified Framework for Factuality Evaluation of LLMs. *arXiv (Cornell University)*.
<https://doi.org/10.48550/arxiv.2405.05583>
- Yu, T., Sonish Sivarajkumar, Kapoor, S., Stolyar, A. V., Polanska, K., McCarthy, K. R., Osterhoudt, H., Wu, X., Shyam Visweswaran, Fu, S., Mathur, P., Cacciamani, G. E., Sun, C., Peng, Y., & Wang, Y. (2024). A framework for human evaluation of large language

models in healthcare derived from literature review. *Npj Digital Medicine*, 7(1).

<https://doi.org/10.1038/s41746-024-01258-7>

Zhang, T., Kishore, V., Wu, F. F., Weinberger, K. Q., & Yoav Artzi. (2019). *BERTScore:*

Evaluating Text Generation with BERT. <https://doi.org/10.48550/arxiv.1904.09675>

Zhao, H., Liu, Z., Wu, Z., Li, Y., Yang, T., Shu, P., Xu, S., Dai, H., Zhao, L., Mai, G., Liu, N., &

Liu, T. (2024). Revolutionizing Finance with LLMs: An Overview of Applications and

Insights. *arXiv (Cornell University)*. <https://doi.org/10.48550/arxiv.2401.11641>