

# Semantic Entropy for Detecting Bias in LLM Responses

Diya Vinod

April 15, 2025

## 1 Introduction

This project proposal focuses on using semantic entropy as a method to investigate bias in Large Language Models (LLMs). One significant challenge that LLMs face is the occurrence of hallucinations. Hallucinations refer to instances where the model generates false, misleading, or nonsensical information that is not based on real-world facts or its training data [1]. Unlike simple factual errors, hallucinations can be highly plausible yet completely fabricated, making them particularly difficult to detect. These inaccuracies not only diminish trust in AI-generated responses but also risk amplifying misinformation, especially in sensitive domains such as medicine, law, and finance [2].

Hallucinations often arise due to uncertainty in the model’s internal knowledge representation. When an LLM encounters a low-confidence or ambiguous query, it does not explicitly recognize gaps in its knowledge. Instead, it interpolates or fabricates information based on patterns learned during training. This process can lead to responses that appear confident but lack factual grounding. The level of uncertainty in a model’s output can vary significantly based on prompt structure, dataset biases, and response generation techniques.

Semantic entropy offers a cost-effective and reliable means of quantifying this uncertainty in LLM-generated text. Unlike traditional confidence scores, which only measure token-level probabilities, semantic entropy captures variation across multiple outputs, allowing us to systematically detect confabulations—instances where the model produces responses with high confidence despite their inaccuracy. By leveraging entropy-based probes, we can assess whether biased or controversial prompts lead to greater semantic variability, providing insights into how LLMs handle uncertainty and bias simultaneously [1, 3].

The research question we attempt to answer is: **“Does the entropy of responses increase for prompts that are more likely to elicit biased answers?”** We anticipate a positive relationship between controversial prompts and semantic entropy. As the question becomes more open-ended and more likely to prompt biased responses, we predict the entropy will also increase.

## 2 Data Collection

Much of the feedback I received on my initial proposal was that it was unclear how prompts would be classified as fact-based vs controversial. I needed to figure out a way to classify these prompts so there was a clear distinction between noncontroversial and controversial prompts.

One of the papers suggested by my peers constructs a dataset that contains controversial questions called DELPHI. This dataset expands on the Quora Question Pairs Dataset. I used these questions to determine the controversial and noncontroversial questions to ask the LLM. By using this dataset, I have used the methods used in the paper to classify my prompts.

Below is the code I used to filter out the most controversial questions vs the least controversial questions. The DELPHI dataset consists of *five* columns.

1. "qid: controlversial question id list corresponding to the original kaggle training dataset"
2. "r1: human annotated strong emotional reaction score (1 as least, 5 as highest)"
3. "r2: human annotated diverse and opposing opinions score (1 as least, 5 as highest)"
4. "Controversial question: Boolean label indicating human annotated result"

Questions that received average ratings of at least 4.0 in both r1 and r2 were labeled "True" in the *controversial question* column. In the code below we get a dataset of just the controversial questions by filtering out rows that have "False" in the *controversial question* column. We also get a dataset of just the noncontroversial questions by filtering out the rows that have values greater than 1.0 in either the "r1" or "r2" column. [4, 5]

```
import pandas as pd
import zipfile
controversial = pd.read_csv('controversial_questions_annotated_id_removed.tsv', sep
↪ ='\t')
controversial.to_csv('controversial_questions_annotated_id_removed.csv', index=False)
controversial.head()
zf = zipfile.ZipFile('questions.csv.zip')
df = pd.read_csv(zf.open('questions.csv'))
df.head()
controversial_true = controversial[controversial['Controversial question'] == True]
merged_qid1 = pd.merge(controversial_true, df, left_on='qid', right_on='qid1',
↪ how='inner')
result_qid1 = merged_qid1[['qid', 'question1']].rename(columns={'question1':
↪ 'question'})
merged_qid2 = pd.merge(controversial_true, df, left_on='qid', right_on='qid2',
↪ how='inner')
result_qid2 = merged_qid2[['qid', 'question2']].rename(columns={'question2':
↪ 'question'})
final_result = pd.concat([result_qid1, result_qid2], ignore_index=True)
final_result.drop_duplicates(inplace=True)
final_result.sort_values(by='qid', inplace=True)
final_result.to_csv('controversial_questions.csv', index=False)
least_controversial = controversial[(controversial['r1'] <= 1.0)&(controversial['r2']
↪ <= 1.0)]
merged_qid1 = pd.merge(least_controversial, df, left_on='qid', right_on='qid1',
↪ how='inner')
result_qid1 = merged_qid1[['qid', 'question1']].rename(columns={'question1':
↪ 'question'})
merged_qid2 = pd.merge(least_controversial, df, left_on='qid', right_on='qid2',
↪ how='inner')
result_qid2 = merged_qid2[['qid', 'question2']].rename(columns={'question2':
↪ 'question'})
final_result = pd.concat([result_qid1, result_qid2], ignore_index=True)
final_result.drop_duplicates(inplace=True)
final_result.sort_values(by='qid', inplace=True)
final_result.to_csv('least_controversial_questions.csv', index=False)
```

After getting the separating the dataset into most controversial questions and least controversial questions we select a few questions from each for our research.

### Low Semantic Entropy Response (Confident, Factual)

These prompts should have low semantic entropy because LLMs are certain about these answers. They are the least controversial questions.

QID	Question
53879	What are some nice 30 days challenge ideas?
220532	Where are the best places to eat in Venice, Italy?
224352	Is eating apple seeds bad for you?
480695	How did Argentina gain independence?
511425	What is engine coolant fluid?

### High Semantic Entropy Response (Uncertain, Hallucinating)

These prompts push for LLMs to provide responses that it may be uncertain about. Using this uncertainty, we can detect hallucinations. The following are questions that may elicit biased responses:

QID	Question
1328	What happens if a dictatorship continues in the present day?
32568	Do you think uneducated people should be allowed to vote?
44601	Is Atticus Finch a good parent?
153814	What is your view/opinion about Brexit?
514520	Are the best things in life truly free?

## 3 Methods

To test our hypothesis that prompts more likely to elicit biased answers have higher entropy, we will employ a comparative analysis using multiple LLMs, including GPT-4o, DeepSeek, and Llama 3.3. We will use the data we collected explained above. With each of the responses, we will calculate the semantic entropy.

### 3.1 Calculating Semantic Entropy in LLM Responses

Below is the equation used to calculate the entropy at the token level, using the probability distribution of tokens in the LLM’s response.

$$H(X) = - \sum_i P(x_i) \log P(x_i) \quad (1)$$

- $H(X)$  is the entropy of the response
- $P(x_i)$  is the probability of token  $x_i$

The sum runs over all tokens in the response.

Steps to compute:

1. Enable probability outputs in the LLM (e.g., using OpenAI’s logprobs and equivalent in other models).

2. For a given prompt, extract the probability distribution over tokens in the generated response.
3. Compute Shannon entropy using the formula above.
4. Compare entropy across different prompts (e.g., fact-based vs. biased prompts).

We will compare semantic entropy across the different prompts and different LLMs. In accordance with our hypothesis we expect that fact-based prompts will have low entropy vs more biased prompts will have the high entropy due to model uncertainty.

## 4 Expected Outcome

We anticipate that prompts categorized as fact-based will exhibit low semantic entropy, as LLMs tend to generate consistent and confident responses when dealing with well-established knowledge. In contrast, controversial or opinion-based prompts are expected to show higher entropy due to the model's uncertainty in navigating subjective or divisive topics. For example, a question like *"Should the United States have free healthcare?"* may lead to varied responses depending on framing, training data, and biases within the model.

However, if our hypothesis is false, we may find that semantic entropy does not strongly correlate with biased or controversial responses. Instead, LLMs might exhibit low entropy even when producing biased responses, indicating that the model is overconfident in its biased outputs rather than uncertain. Alternatively, semantic entropy might fluctuate based on factors unrelated to bias, such as prompt phrasing, model architecture, or fine-tuning

## 5 Minimum Viable Example

For a minimum viable example, I will be doing a simple test of my methodology in GPT-4o. I will compute the semantic entropy of two fact-based prompts and two more biased prompts. Then I will compare the two outputs to determine if there are significant differences.

## 6 Timeline

1. *Week 1: Literature Review and Finalize Methods* - Read about previous works in semantic uncertainty. Understand the math and theory behind the semantic entropy equation.
2. *Week 2: Data Collection* - Formulate fact-based questions and controversial questions to ask LLMs. Obtain access to the various LLMs and ask the questions.
3. *Week 3 & Week 4: Data Analysis* - Calculate the semantic entropy of the responses and do comparative analysis.
4. *Week 5: Conclusion and Report* - Summarize process and key findings in a paper.

## 7 References

### References

- [1] Farquhar, S., Kossen, J., Kuhn, L., & Gal, Y. (2024). Detecting hallucinations in large language models using semantic entropy. *Nature*, 630, 625–630. <https://doi.org/10.1038/s41586-024-07421-0>
- [2] IBM. (2023, September 1). *AI hallucinations*. Ibm.com. <https://www.ibm.com/think/topics/ai-hallucinations>
- [3] Kossen, J., Han, J., Razzak, M., Schut, L., Malik, S., & Gal, Y. (2024). Semantic Entropy Probes: Robust and Cheap Hallucination Detection in LLMs. *arXiv*. <https://arxiv.org/abs/2406.15927>
- [4] Sun, D. Q., Abzaliev, A., Kotek, H., Xiu, Z., Klein, C., & Williams, J. D. (2023). DELPHI: Data for Evaluating LLMs’ Performance in Handling Controversial Issues. *arXiv preprint arXiv:2310.18130*. <https://arxiv.org/abs/2310.18130>
- [5] Farquhar, S., Kossen, J., Kuhn, L., & Gal, Y. (2024, June 19). *Detecting hallucinations in large language models using semantic entropy*. Oxford Applied and Theoretical Machine Learning Group. [https://oatml.cs.ox.ac.uk/blog/2024/06/19/detecting\\_hallucinations\\_2024.html](https://oatml.cs.ox.ac.uk/blog/2024/06/19/detecting_hallucinations_2024.html)

## 8 LLM Usage

I knew I wanted to pursue the topic of hallucinations after reading the paper, “Detecting hallucinations in large language models using semantic entropy.” Once I knew I wanted to focus on this topic, I had ChatGPT generate ideas for ways I could expand upon the idea. I used ChatGPT to help me gain a better understanding of the concept of semantic entropy. I also used it to provide equations and validate my methods.