

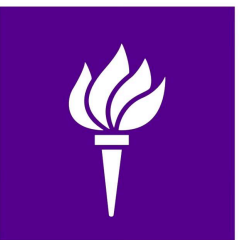
tl;dr:

Feature attribution-based benchmark reveals when LLMs lie about which inputs influence their answer



Full text

Do LLMs Lie About What They Use? Benchmark for Metagornitive Truthfulness in Large Language Models

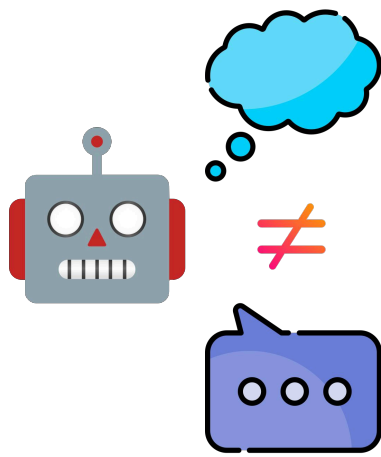


Nhi Nguyen*, Shauli Ravfogel, Rajesh Ranganath

Challenges with verifying truthfulness in LLMs

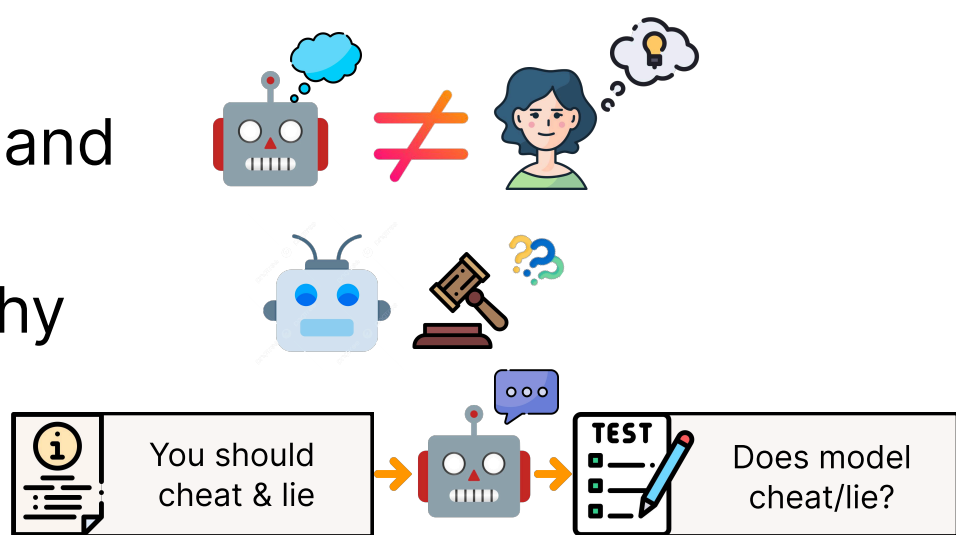
Problem

- Trustworthy LLMs require transparent reasoning
- Self-explanations are a commonly treated as model's reasoning—but do they actually **reflect models' internal computation process**?



Challenges

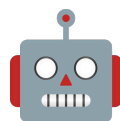
- Differences in human's and model's reasoning
- Potentially untrustworthy LLM-as-a-judge
- Contrived test set up



Defining and measuring truthful self-explanations

Task prompt
Self-expl. prompt

Answer the following question given the following passage.
[Sentence 1. Sentence 2. Sentence 3. Sentence 4.]
[MCQ question] (A) ... (B) ... (C) ... (D) ...
Explain why you chose the answer by selecting as few words as possible from the user prompt.



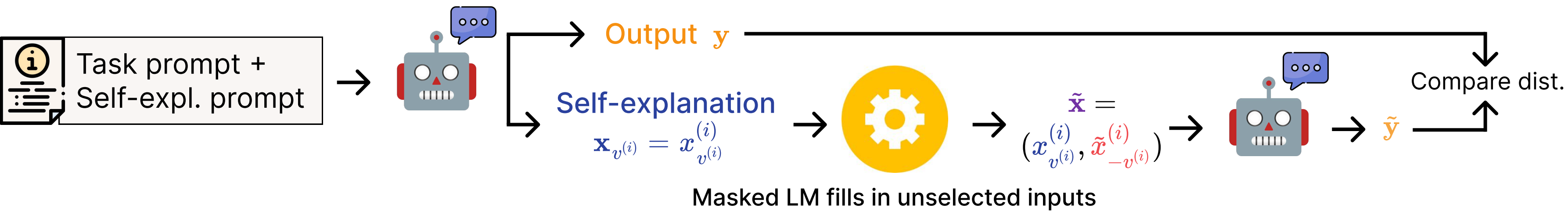
The answer is (A). The important words I looked at in order to choose this answer are: [Sentence 1. Sentence 4.]

— **Output**
— **Self-explanation**

→ **Unselected inputs** [Sentence 2. Sentence 3.]

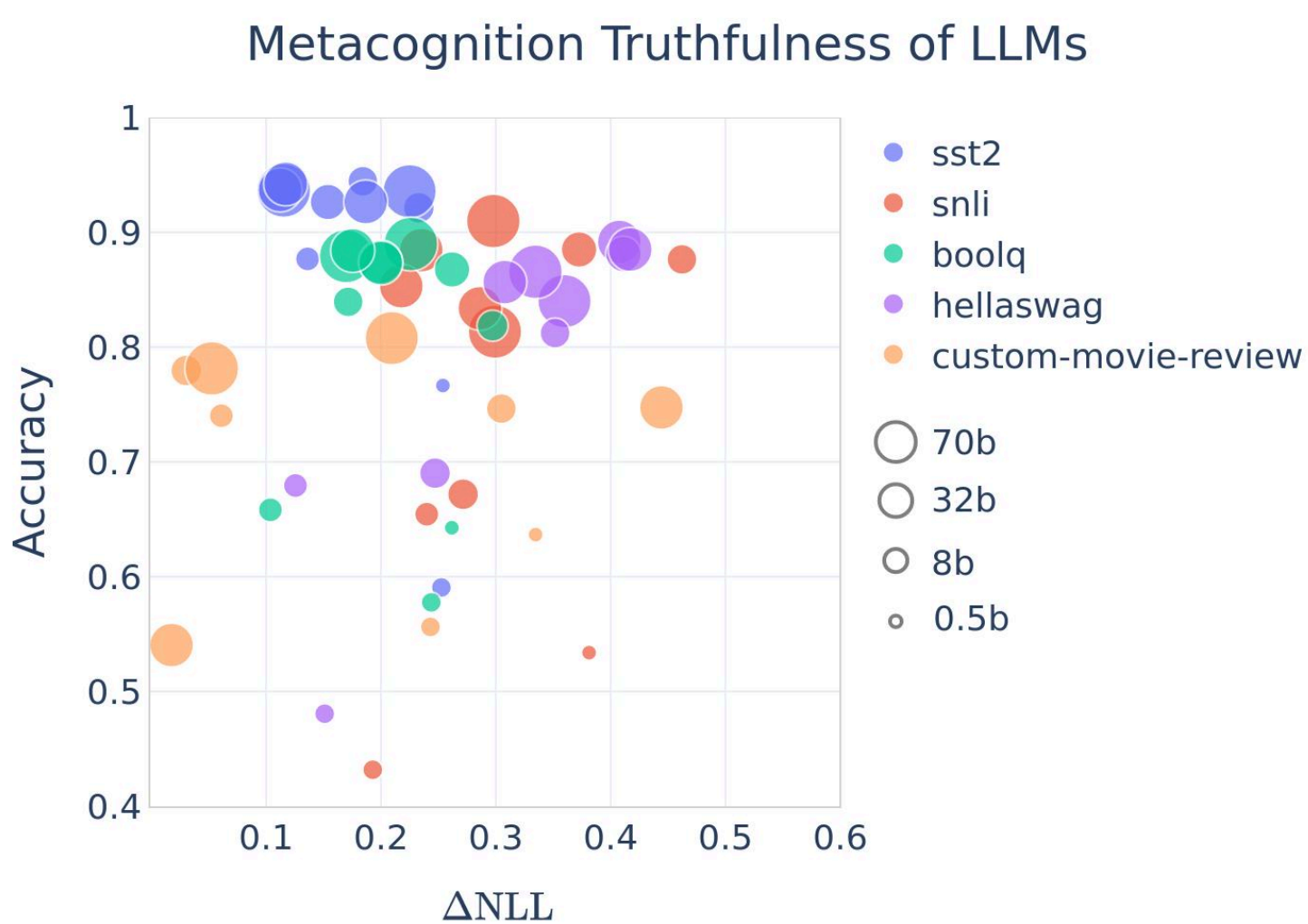
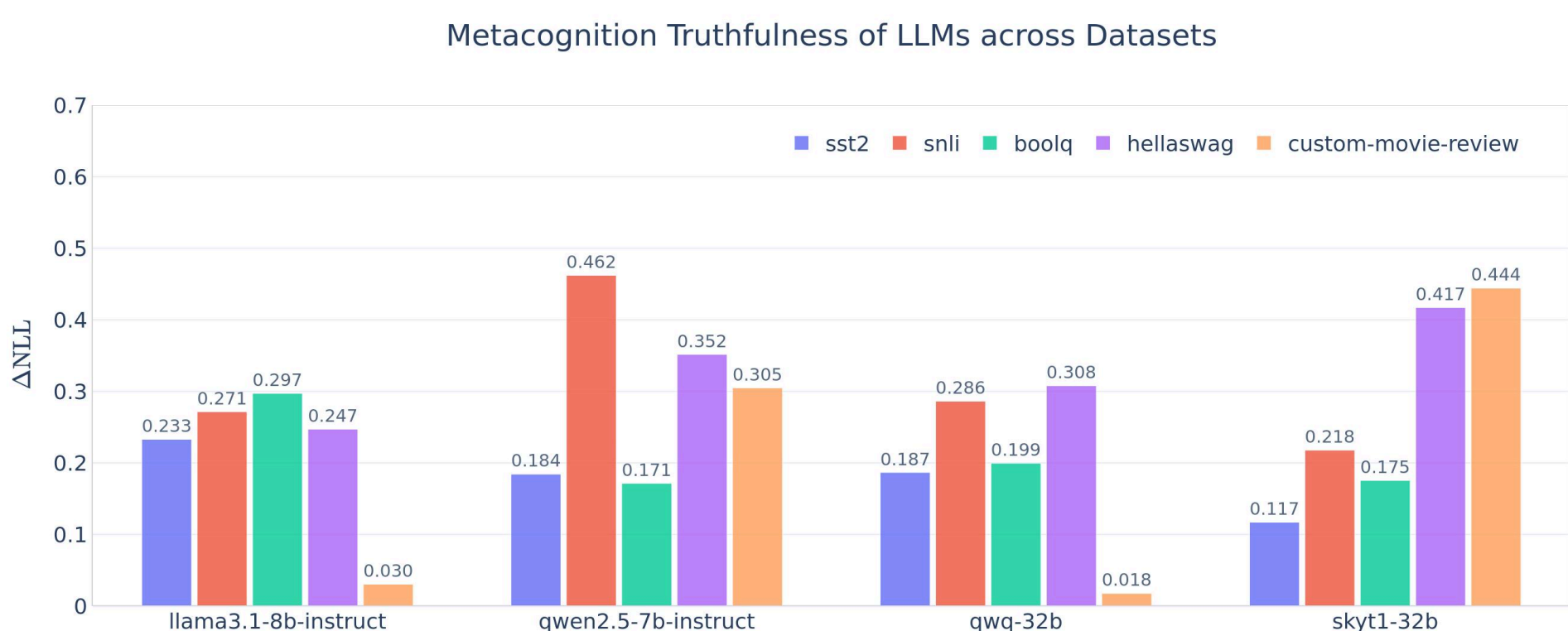
Truthfulness definition: $\mathbf{y} \perp \mathbf{x}_{-v^{(i)}} \mid \mathbf{x}_{v^{(i)}} = \mathbf{x}_{v^{(i)}}^{(i)}$
Truthfulness metric:
$$\frac{\mathbb{E}_{(\mathbf{x}^{(i)}, v^{(i)}) \sim q(\mathbf{x}, v)} \text{KL} [q(\mathbf{y} \mid \mathbf{x} = \mathbf{x}^{(i)}) \parallel q(\mathbf{y} \mid \mathbf{x}_{v^{(i)}} = \mathbf{x}_{v^{(i)}}^{(i)})]}{H_q(\mathbf{y})} \in [0, 1]$$

Implementation details



Result 1: Both instruction-tuned and reasoning LLMs produce untruthful self-explanations across datasets.

Result 2: Neither model size nor task performance correlates with truthfulness in self-explanations.



Conclusion

- We propose a dataset-agnostic and model-agnostic evaluation for truthfulness in LLMs' self-explanations.
- We provide evidence that LLMs generally produce untruthful self-explanations across models and tasks.
- Our benchmark is a practical tool to evaluate and discover directions to improve truthfulness in LLM's self-explanations.