SaGE: Quantifying Moral Consistency in Large Language Models

Vamshi Krishna Bonagiri, Sreeram Vennam, Priyanshul Govil, Manas Gaur, Ponnurangam Kumaraguru







HYDERABAD

Large Language Models are Inconsistent



Large Language Models are Inconsistent

Semantic Consistency: the ability to make consistent decisions in semantically equivalent contexts. i.e , Semantically equivalent questions should yield semantically equivalent answers

Elazar et al. 2021



Claim: LLMs do not have this ability, and can give contradictory answers to paraphrased questions

Why is it a problem?



Behave in unexpected ways - hindering trust and reliability (Safety Risk)

Problem Statement

How can we reliably measure LLM Consistency in free text generation?

Methodology ► Pipeline



Semantic graph

SaGE Score

Methodology ► Rules of Thumb



An example from the Moral Integrity Corpus

Methodology > Paraphrase and RoT Generation

Generating Paraphrases

- Few shot prompting Vicuna-13B
- Filtering out good quality paraphrases using ParaScore (> 0.8)
- Semantic Similarity: Do the paraphrases mean the same thing?
- Lexical Divergence: Do they have enough variation in vocabulary?

Generating RoTs

- Few shot prompting with proper instructions seemed to produce good enough results.

Our prompt: RoT generation

A Rule of Thumb (RoT) is a fundamental judgment about right and wrong behavior. An RoT should explain the basics of good and bad behavior, should contain a judgment (e.g. "you should") and an action (e.g. "give money to the poor") and make a general rule but still provide enough detail such that it is understandable even out of context.

Instruction: Generate the rule of thumb for the task question-answer pair: be as concise as possible.

Question: When is it worth it to fight a losing battle?

Answer: It is never worth it.

Rule of Thumb: Cowardice saves you from fights you cannot win.

. . .

Question: <question> Answer: <answer> Rule of Thumb:

Methodology ► pipeline



10K Questions, 5 Paraphrases each

Methodology ► Semantic Graph

Semantic Graphs:

- Textual responses are converted into semantic embeddings using SBERT DeBERTa, fine tuned on NLI datasets.
- Each sentence representation is a node in the graph
- Distance between two nodes is the cosine distance between their semantic embeddings



f

Methodology > Semantic Graph Entropy

Graph Entropy as a measure for consistency

- Less entropy = Consistent
- More entropy = Inconsistent

f: Information functional

$$v_i) = \sum_{i=1}^{n} \sin(v_i, v_i)$$

$$H(p) = -\sum_{i=1}^{n} p_i \log(p_i)$$

$$p(v_i) = \frac{f(v_i)}{\sum_{j=1}^{|V|} f(v_j)},$$

Probability function

$$\operatorname{SaGE}(G_s) = 1 - \frac{I(G_s)}{\log n}$$

Dehmer and Mowshowitzt., 2011

Results ► LLMs are inconsistent

Model	BLEU		ROUGE		BERTScore		SaGE	
	Ans	RoT	Ans	RoT	Ans	RoT	Ans	RoT
opt-125m	0.011	0.012	0.138	0.127	0.355	0.352	0.243	0.252
opt-1.3b	0.009	0.010	0.133	0.119	0.369	0.362	0.263	0.268
opt-2.7b	0.008	0.011	0.135	0.127	0.382	0.378	0.277	0.284
opt-6.7b	0.007	0.012	0.130	0.129	0.385	0.382	0.282	0.290
opt-13b	0.008	0.012	0.139	0.135	0.412	0.408	0.312	0.318
Mistral-7B-Instruct-v0.1	0.016	0.015	0.151	0.150	0.499	0.493	0.405	0.407
falcon-7b-instruct	0.027	0.016	0.194	0.159	0.648	0.621	0.584	0.563
Llama-2-7b-chat-hf	0.073	0.020	0.296	0.170	0.564	0.546	0.362	0.452
Llama-2-13b-chat-hf	0.084	0.020	0.261	0.176	0.660	0.635	0.595	0.575
GPT-3.5 Turbo †	0.056	0.015	0.217	0.151	0.613	0.529	0.681	0.478
GPT-4 †	0.055	0.0172	0.246	0.166	0.568	0.486	0.641	0.438

SOTA LLMs have consistency around 0.575

Results > Human Annotations

Metric	Answers	RoTs	
BLEU	0.391	0.412	
ROUGE	0.459	0.476	
BERTScore	0.522	0.527	
SaGE	0.561	0.592	

- Human annotations for 500 data points
- 0 and 1 for consistent/inconsistent for each pair
- Average of these annotations compared with Metric scores

Results Is Consistency dependent on temperature?



Probably not

Results - Generalization



(Commonsense Reasoning)

Model	Trut	hfulQA	HellaSwag		
	SaGE	Accuracy	SaGE	Accuracy	
opt-125m	0.258	0.357	0.164	0.313	
opt-1.3b	0.258	0.260	0.162	0.537	
opt-2.7b	0.282	0.374	0.151	0.614	
opt-6.7b	0.285	0.351	0.156	0.687	
opt-13b	0.315	0.341	0.146	0.712	
Mistral-7B	0.421	0.567	0.529	0.756	
falcon-7b	0.577	0.343	0.289	0.781	
Llama-2-7b	0.452	0.388	0.563	0.786	
Llama-2-13b	0.559	0.374	0.520	0.819	

Surprisingly, Accuracy on benchmarks does not correlate with consistency!

Results ► Improvement?

Model	BLEU	ROUGE	BERT Score	SaGE
GPT-3.5	0.015	0.151	0.529	0.438
GPT-3.5 with RoT prompt- ing	0.018	0.169	0.565	0.548

Our prompt: RoT-based answer generation

Instruction: Answer the following question.
Keep in mind this rule of thumb, <RoT>
Question: <question>

Answer:

RoT prompting shows improvement as expected

Conclusion and Future Work

- We introduce the Moral Consistency Corpus, with 50K moral questions and LLM responses to them
- We introduce the SaGE framework for measuring consistency of an LLM
- Results show that
 - Current LLMs are inconsistent
 - Consistency and Accuracy are not the same problem, and this problem has to be investigated further. Models need to be evaluated for consistency separately.
 - Providing simple rules to follow increases consistency, hinting at the potential for methods such as Retrieval Augmented Generation (RAG) for improvement in consistency.

Thank You!

Questions? Please feel free to reach out to



vamshi.b@research.iiit.ac.in



@VictorKnox99





Paper



