

Logical Consistency of Large Language Models in Fact-Checking

Bishwamitra Ghosh, Sarah Hasan, Naheed Anjum Arafat, Arijit Khan



MAX PLANCK INSTITUTE
FOR SOFTWARE SYSTEMS



AALBORG
UNIVERSITY

LLMs are impressive at general language understanding, yet they suffer from inconsistency

LLMs are impressive at general language understanding, yet they suffer from inconsistency

What is consistency?

LLMs are impressive at general language understanding, yet they suffer from inconsistency

What is consistency?

Existing Works: **Similar response to semantically similar prompt**

LLMs are impressive at general language understanding, yet they suffer from inconsistency

What is consistency?

Existing Works: **Similar response to semantically similar prompt**

Paraphrasing

Berlin is the capital of Germany

Germany's capital is Berlin

$\text{LLM}(\text{Berlin is the capital of Germany}) = \text{LLM}(\text{Germany's capital is Berlin})$

What is Logical Consistency?

Our Proposal

Response is consistent with logical changes of the prompt

- ▶ Similar response to logically equivalent prompt
- ▶ Different response to logically different prompt
- ▶ Response should adhere to formal logic

What is Logical Consistency?

Our Proposal

Response is consistent with logical changes of the prompt

- ▶ Similar response to logically equivalent prompt
- ▶ Different response to logically different prompt
- ▶ Response should adhere to formal logic

Negation

Berlin is the capital of Germany

Berlin is **not** the capital of Germany

$\text{LLM}(\text{Berlin is the capital of Germany}) \neq \text{LLM}(\text{Berlin is **not** the capital of Germany})$

Conjunction

Berlin is the capital of Germany **and** US embassy is in Berlin

Berlin is the capital of Germany

US embassy is in Berlin

Conjunction

Berlin is the capital of Germany **and** US embassy is in Berlin

Berlin is the capital of Germany

US embassy is in Berlin

$$\text{LLM}\left(\begin{array}{l} \text{Berlin is the capital of Germany} \\ \text{and} \\ \text{US embassy is in Berlin} \end{array}\right) = \begin{array}{l} \text{LLM}(\text{Berlin is the capital of Germany}) \\ \text{and} \\ \text{LLM}(\text{US embassy is in Berlin}) \end{array}$$

- ▶ Logical consistency on complex logical queries with negation, conjunction, and disjunction operators
- ▶ As a specific test bed, we consider the task of [fact-checking in knowledge graphs](#) (KGs) using LLMs

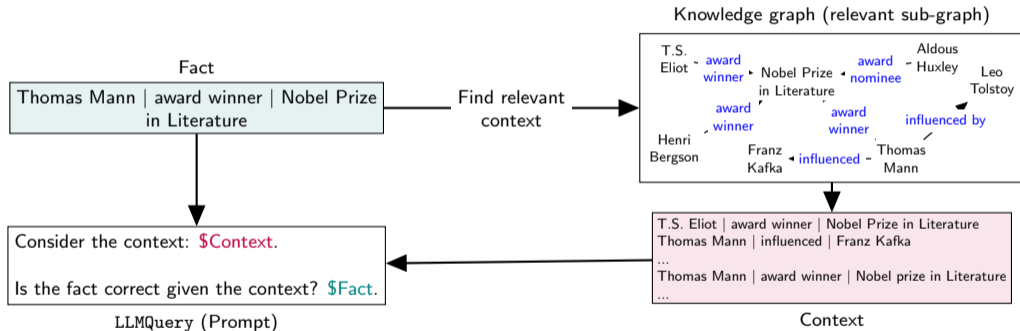
- ▶ Logical consistency on complex logical queries with negation, conjunction, and disjunction operators
- ▶ As a specific test bed, we consider the task of [fact-checking in knowledge graphs](#) (KGs) using LLMs

Benchmark

Assessment

Improvement

Our Framework: LLM in fact-checking with KG



Primitive operators

$$\text{LLM}(\neg p) = \neg \text{LLM}(p)$$

$$\text{LLM}(p \vee q) = \text{LLM}(p) \vee \text{LLM}(q)$$

$$\text{LLM}(p \wedge q) = \text{LLM}(p) \wedge \text{LLM}(q)$$

Primitive operators

$$\text{LLM}(\neg p) = \neg \text{LLM}(p)$$

$$\text{LLM}(p \vee q) = \text{LLM}(p) \vee \text{LLM}(q)$$

$$\text{LLM}(p \wedge q) = \text{LLM}(p) \wedge \text{LLM}(q)$$

Disjunctive normal form (DNF): A DNF fact $q = \bigvee_{i=1}^n c_i$, where $c_i = \bigwedge_{j=1}^{i_m} e_{ij}$

$$\text{LLM}(q) = \bigvee_{i=1}^n \left(\bigwedge_{j=1}^{i_m} \text{LLM}(e_{ij}) \right)$$

Commutative law

$$\text{LLM}(p \vee q) = \text{LLM}(q \vee p)$$

$$\text{LLM}(p \wedge q) = \text{LLM}(q \wedge p)$$

Associative law

$$\text{LLM}((p \vee q) \vee s) = \text{LLM}(p \vee (q \vee s))$$

$$\text{LLM}((p \wedge q) \wedge s) = \text{LLM}(p \wedge (q \wedge s))$$

Distributive law

$$\text{LLM}(p \wedge (q \vee s)) = \text{LLM}((p \wedge q) \vee (p \wedge s))$$

$$\text{LLM}(p \vee (q \wedge s)) = \text{LLM}((p \vee q) \wedge (p \vee s))$$

... De-Morgan's Laws and First-order logic.

Model	Dataset	Fact	Accuracy		Logical Consistency	
			Before FT ¹	After FT	Before FT	After FT
Llama2-13B	FreebaseLFC	$p, \neg p$	0.90		0.81	
		$p \wedge q$	0.61		0.67	
		$p \vee q$	0.73		0.73	
	NELLFC	$p, \neg p$	0.88		0.76	
		$p \wedge q$	0.38		0.69	
		$p \vee q$	0.73		0.73	
	WikiLFC	$p, \neg p$	0.96		0.92	

¹FT = Fine-tuning

Model	Dataset	Fact	Accuracy		Logical Consistency	
			Before FT ¹	After FT	Before FT	After FT
Llama2-13B	FreebaseLFC	$p, \neg p$	0.90		0.81	
		$p \wedge q$	0.61		0.67	
		$p \vee q$	0.73		0.73	
	NELLFC	$p, \neg p$	0.88		0.76	
		$p \wedge q$	0.38		0.69	
		$p \vee q$	0.73		0.73	
	WikiLFC	$p, \neg p$	0.96		0.92	

¹FT = Fine-tuning

- ▶ An LLM is consistent on a simple atomic fact if it is accurate both on the fact and its negation
- ▶ For a complex DNF fact, the LLM is consistent if it is accurate on the DNF fact as well as on all constituent atomic facts

Model	Dataset	Fact	Accuracy		Logical Consistency	
			Before FT	After FT	Before FT	After FT
Llama2-13B	FreebaseLFC	$p, \neg p$	0.90	0.93	0.81	0.86
		$p \wedge q$	0.61	0.93	0.67	0.83
		$p \vee q$	0.73	0.76	0.73	0.97
	NELLFC	$p, \neg p$	0.88	0.97	0.76	0.93
		$p \wedge q$	0.38	0.89	0.69	0.88
		$p \vee q$	0.73	0.76	0.73	0.94
	WikiLFC	$p, \neg p$	0.96	0.96	0.92	0.93

- ▶ **Assessment.** LLMs are not always logically consistent with their generation – consistency decreases as the query complexity increases.

- ▶ **Assessment.** LLMs are not always logically consistent with their generation – consistency decreases as the query complexity increases.
- ▶ **Improvement.** Instruction prompting is not sufficient to improve logical consistency in LLMs: smaller models require instruction fine-tuning, while larger models may suffice with instruction prompting.

- ▶ **Assessment.** LLMs are not always logically consistent with their generation – consistency decreases as the query complexity increases.
- ▶ **Improvement.** Instruction prompting is not sufficient to improve logical consistency in LLMs: smaller models require instruction fine-tuning, while larger models may suffice with instruction prompting.
 - ▶ **Generalization.** Fine-tuning for logical consistency in one dataset can generalize to other datasets and queries with more logical operators.

- ▶ **Assessment.** LLMs are not always logically consistent with their generation – consistency decreases as the query complexity increases.
- ▶ **Improvement.** Instruction prompting is not sufficient to improve logical consistency in LLMs: smaller models require instruction fine-tuning, while larger models may suffice with instruction prompting.
 - ▶ **Generalization.** Fine-tuning for logical consistency in one dataset can generalize to other datasets and queries with more logical operators.
- ▶ **Benchmark.** Fact-checking in KG provides a flexible benchmark to test LLMs on logical queries of varying complexity.

- ▶ **Assessment.** LLMs are not always logically consistent with their generation – consistency decreases as the query complexity increases.
- ▶ **Improvement.** Instruction prompting is not sufficient to improve logical consistency in LLMs: smaller models require instruction fine-tuning, while larger models may suffice with instruction prompting.
 - ▶ **Generalization.** Fine-tuning for logical consistency in one dataset can generalize to other datasets and queries with more logical operators.
- ▶ **Benchmark.** Fact-checking in KG provides a flexible benchmark to test LLMs on logical queries of varying complexity.

- ▶ Logical inconsistency is a critical issue for LLMs despite their impressive language understanding ability
- ▶ Propose a framework to assess the logical consistency of LLMs on complex fact-check queries from KGs
- ▶ Demonstrate how supervised fine-tuning can improve the logical consistency of LLMs



Paper