

# **RAG-HAT: A Hallucination-Aware Tuning Pipeline for LLM in Retrieval-Augmented Generation**

*Juntong Song, Xingguang Wang, Juno Zhu, Yuanhao Wu, Xuxin Cheng, Randy Zhong, Cheng Niu*

# RAGTruth: A Hallucination Corpus for Developing Trustworthy Retrieval-Augmented Language Models

Cheng Niu<sup>1</sup>, Yuanhao Wu<sup>1</sup>, Juno Zhu<sup>1</sup>, Siliang Xu<sup>1</sup>, Kashun Shum<sup>1</sup>,  
Randy Zhong<sup>1</sup>, Juntong Song<sup>1</sup>, and Tong Zhang<sup>2</sup>

<sup>1</sup>NewsBreak

<sup>2</sup>University of Illinois Urbana-Champaign  
cheng.niu@newsbreak.com

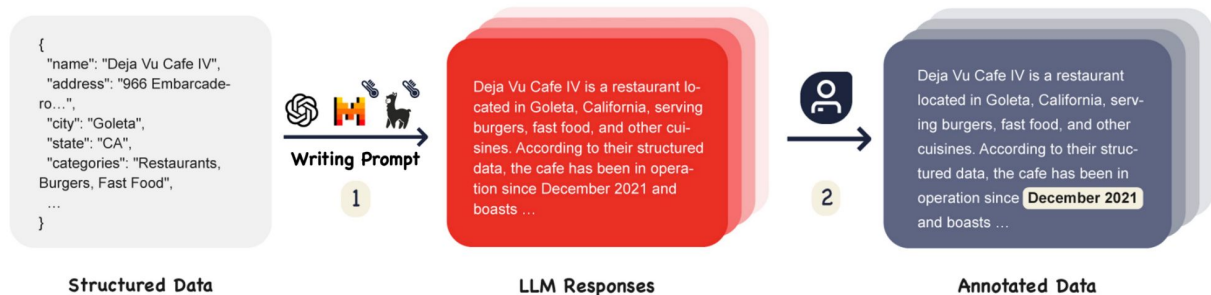


Figure 1: Data gathering pipeline. Taking a data-to-text writing task as an example, our data gathering pipeline includes 2 steps: 1) response generation. We generated responses with multiple LLMs and natural prompts. 2) human annotation. Human labeler annotated hallucinated spans in LLM responses.

# RAGTruth: A Hallucination Corpus for Developing Trustworthy Retrieval-Augmented Language Models



Cheng Niu<sup>1</sup>, Yuanhao Wu<sup>1</sup>, Juno Zhu<sup>1</sup>, Siliang Xu<sup>1</sup>, Kashun Shum<sup>1</sup>,  
Randy Zhong<sup>1</sup>, Juntong Song<sup>1</sup>, and Tong Zhang<sup>2</sup>

<sup>1</sup>NewsBreak

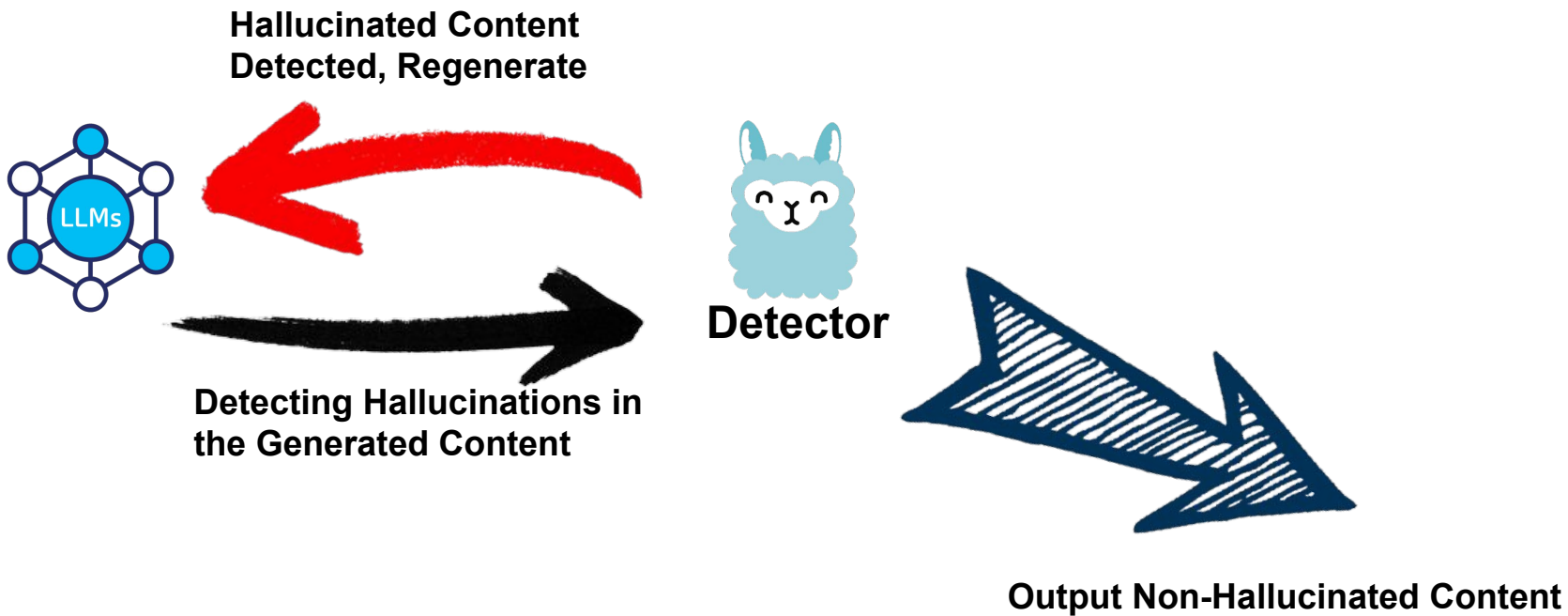
<sup>2</sup>University of Illinois Urbana-Champaign  
cheng.niu@newsbreak.com

Methods	QUESTION ANSWERING			DATA-TO-TEXT WRITING			SUMMARIZATION			OVERALL		
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
Prompt <sub>gpt-3.5-turbo</sub>	18.8	84.4	30.8	65.1	95.5	77.4	23.4	89.2	37.1	37.1	92.3	52.9
Prompt <sub>gpt-4-turbo</sub>	33.2	90.6	45.6	64.3	<b>100.0</b>	78.3	31.5	97.6	47.6	46.9	97.9	63.4
SelfCheckGPT <sub>gpt-3.5-turbo</sub>	35.0	58.0	43.7	68.2	82.8	74.8	31.1	56.5	40.1	49.7	71.9	58.8
LMvLM <sub>gpt-4-turbo</sub>	18.7	76.9	30.1	68.0	76.7	72.1	23.3	81.9	36.2	36.2	77.8	49.4
Finetuned Llama-2-13B	<b>61.6</b>	76.3	<b>68.2</b>	<b>85.4</b>	91.0	<b>88.1</b>	<b>64.0</b>	54.9	<b>59.1</b>	<b>76.9</b>	80.7	<b>78.7</b>

Table 5: The response-level hallucination detection performance for each baseline method across different tasks and different models.

Methods	QUESTION ANSWERING			DATA-TO-TEXT WRITING			SUMMARIZATION			OVERALL		
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
Prompt Baseline <sub>gpt-3.5-turbo</sub>	7.9	25.1	12.1	8.7	45.1	14.6	6.1	33.7	10.3	7.8	35.3	12.8
Prompt Baseline <sub>gpt-4-turbo</sub>	23.7	52.0	32.6	17.9	<b>66.4</b>	28.2	14.7	<b>65.4</b>	24.1	18.4	<b>60.9</b>	28.3
Finetuned Llama-2-13B	<b>55.8</b>	<b>60.8</b>	<b>58.2</b>	<b>56.5</b>	50.7	<b>53.5</b>	<b>52.4</b>	30.8	<b>38.8</b>	<b>55.6</b>	50.2	<b>52.7</b>

Table 6: The span-level detection performance for each baseline method across different tasks and different models.



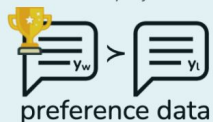
# Hallucinations still persist, even after the model has undergone delicate instruction tuning ...

Model	QUESTION ANSWERING			DATA-TO-TEXT WRITING			SUMMARIZATION			OVERALL	
	# Resp.	# Span	Density	# Resp.	# Span	Density	# Resp.	# Span	Density	# Resp.	# Span
GPT-3.5-turbo-0613	75	89	0.12	272	384	0.18	54	60	0.05	401	533
GPT-4-0613	48	51	0.06	290	354	0.27	74	80	0.08	406	485
Llama-2-7B-chat	510	1010	0.59	888	1775	1.27	434	517	0.58	1832	3302
Llama-2-13B-chat	399	654	0.48	983	2803	1.53	295	342	0.41	1677	3799
Llama-2-70B-chat <sup>†</sup>	320	529	0.40	863	1834	1.15	212	245	0.26	1395	2608
Mistral-7B-Instruct	378	594	0.59	958	2140	1.51	617	828	0.86	1953	3562

Table 3: Hallucination counts and density of models. †: We used 4-bit quantized version of Llama-2-70B-chat.

## Direct Preference Optimization (DPO)

x: "write me a poem about  
the history of jazz"



maximum  
likelihood

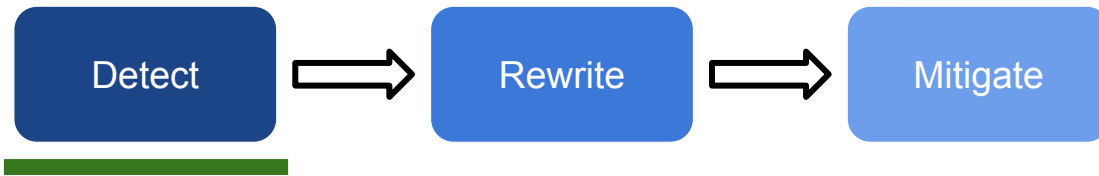


**Reinforcement learning becomes a clear choice when hallucinations still persist, even after the model has undergone delicate instruction tuning.**

In RAG-HAT, We conduct Direct Performance Alignment on the selected LLM. It's a widely used technology that can do preference alignment without explicitly training a reward model.

Our choice is based on the following pain points we identified during our investigations:

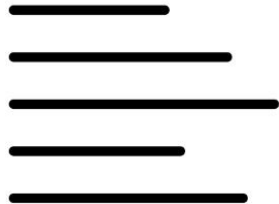
1. First, organizing hallucination mitigation tasks into a preference dataset for reward modeling is challenging, as annotators often struggle with determining which type of hallucination is more severe. In our previous work on RAGTruth, our annotators spent a lot of time struggling to achieve consensus.
2. Second, treating "not hallucinating" as a simple preference is problematic because it essentially requires LLMs to be "always correct," which is an overly rigid expectation. To address this pain point, we will introduce some adjustments later to make the alignment process more natural.



Generated Content



Split into Sentences!!



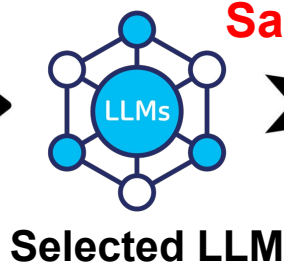
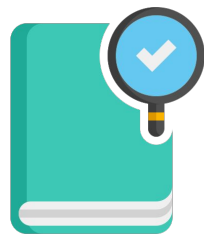
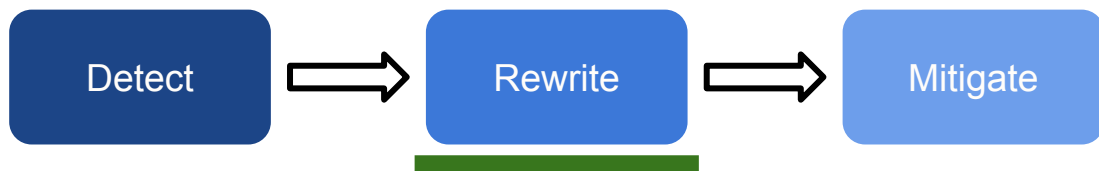
Hallucinated



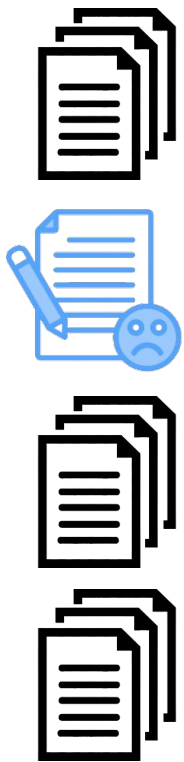
Non-hallucinated



Interpretations



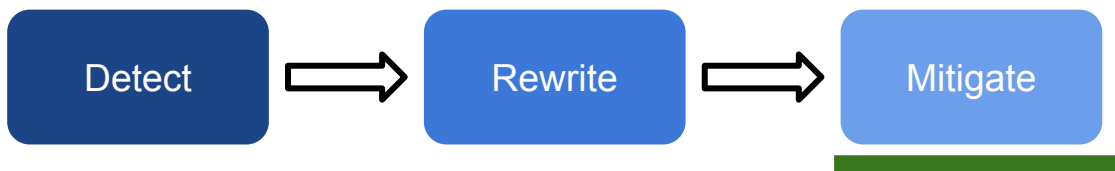
**Sample!**



**Rewrite!**







Preference Dataset



Model Hallucinates Less

## Defensive Advice

... According to their structured data, they offer a variety of beers and a Honey Avocado Ale that is **popular among customers** ...

```
... {  
  "review_stars": 5.0,  
  "review_date": "2021-12-17 02:46:28",  
  "review_text": "Beautiful views, great  
beer, and family friendly. Highly  
recommend the Honey  
Avocado Ale ...
```



Some might argue that the claim of the Honey Avocado Ale being "popular among customers" could be considered a hallucination since the JSON data does not explicitly state its popularity ...

Figure 1: An Example of Defensive Advice: The LLM made a minor extension partially based on the provided references. Defensive advice highlights that the statement is not well supported.

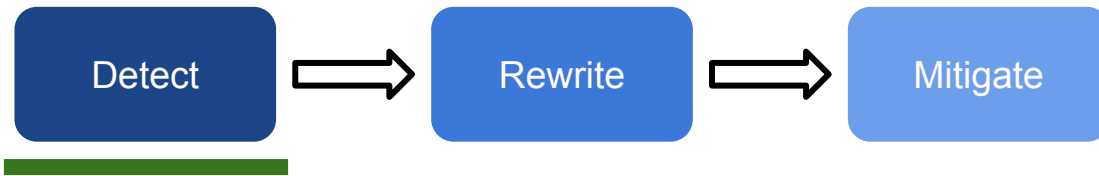
## Overly Cautious Penalization



Teaching yourself can be done by learning something new with each task you complete, as mentioned in passage 1. **However, taking notes is not always necessary.** To learn a language quickly, studying or working abroad can be helpful, as mentioned in passage 2. Wikipedia ...



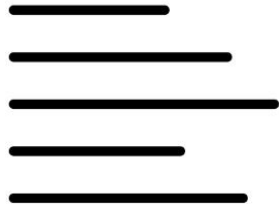
Teaching yourself can be done by learning something new with each task you complete, as mentioned in passage 1. ~~However, taking notes is not always necessary.~~ To learn a language quickly, studying or working abroad can be helpful, as mentioned in passage 2. Wikipedia ...



Generated Content



Split into Sentences!!



Hallucinated



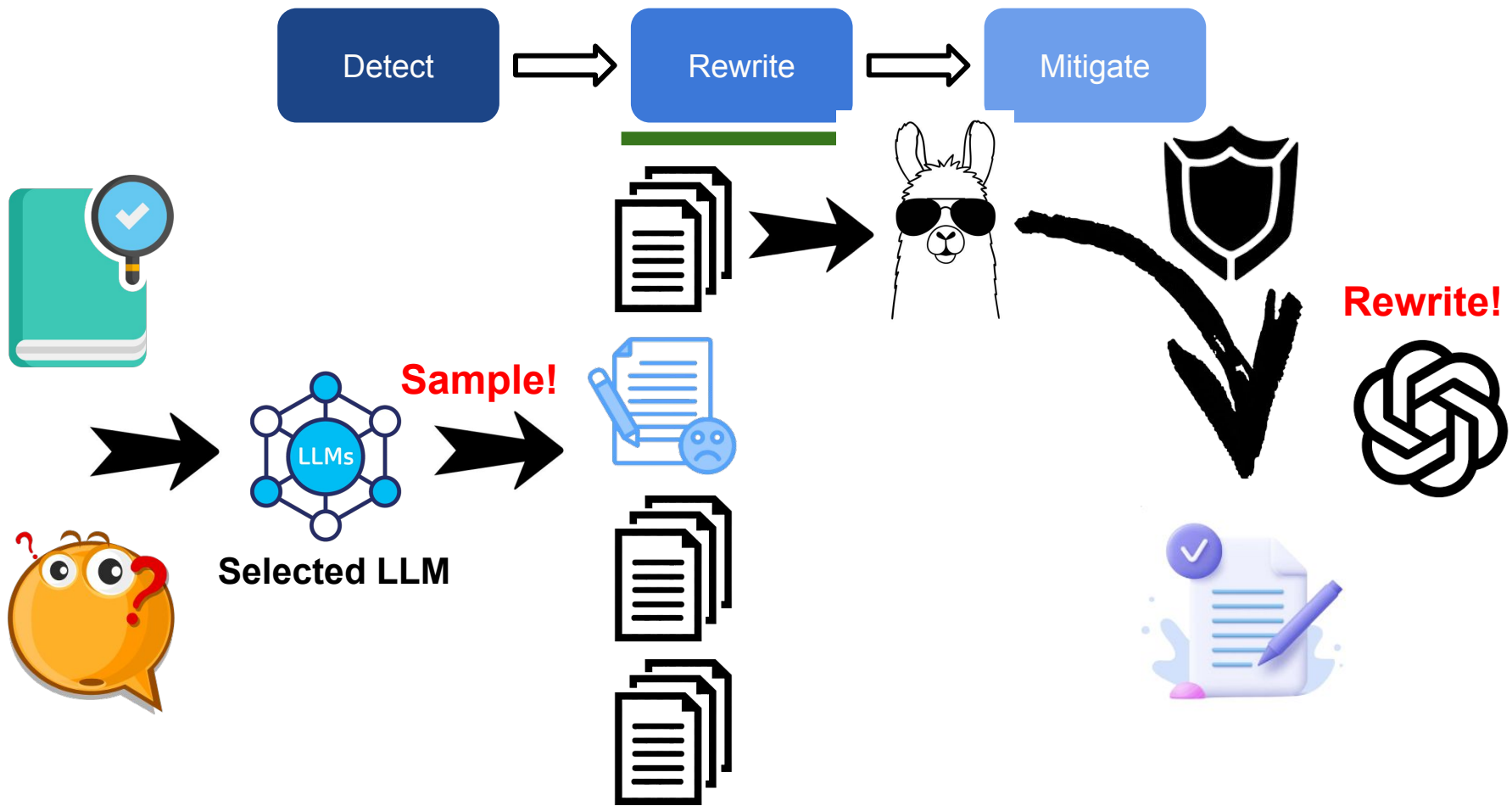
Non-hallucinated



Interpretations



Defensive Advice



DATASET	METHOD	Detector	GPT-4 Turbo	Human	Average
RAGTruth Test Set	Qwen	36.9(-)	51.3(-)	34.4(-)	40.9(-)
	Qwen(Regenerate)	-	44.2(↓13.8%)	-	44.2(↓13.8%)
	RAG-HAT	22.7(↓ <b>38.5%</b> )	41.3(↓ <b>19.5%</b> )	25.7(↓ <b>25.3%</b> )	29.9(↓ <b>26.9%</b> )
WebGLM 1000	Qwen	21.3(-)	46.7(-)	-	34(-)
	Qwen(Regenerate)	-	38.8(↓17.0%)	-	38.8(↓17.0%)
	RAG-HAT	12.0(↓ <b>43.7%</b> )	37.9(↓ <b>19.0%</b> )	-	24.9(↓ <b>26.8%</b> )

Table 3: Hallucination Rate: 1,000-Example WebGLM Set and RAGTruth Test Set (Total 450 Examples): Our detection model cannot fairly benchmark the hallucination rate of the regeneration approach since it serves as the trigger for regeneration.

<b>PAIRED METHOD</b>	<b>WIN RATE</b> (GPT-4 Turbo)
<b>RAG-HAT</b> (full) :: (w/o defensive, w/ OCP)	<b>51.5</b>
<b>RAG-HAT</b> (full) :: (w/o defensive, w/o OCP)	<b>54.1</b>

Table 6: Impact of Training Dataset Composition on Answer Quality: Pairwise Comparison

<b>DATASET</b>	<b>METHOD</b>	<b>GPT-4 Turbo</b>	<b>Human</b>
RAGTruth Dataset	Qwen	41.1	33.2
	<b>RAG-HAT</b>	<b>57.3</b>	<b>40.8</b>
WebGLM 1000	Qwen	39.5	-
	<b>RAG-HAT</b>	<b>58.5</b>	-

Table 5: Answer Quality Win Rates: 1,000-Example WebGLM Set and RAGTruth Test Set

**Thank You!**