

# AutoHall: Automated Hallucination Dataset Generation for Large Language Models

Anonymous ACL submission

## Abstract

Large language models (LLMs) have garnered widespread applications across various domains due to their powerful language understanding and generation capabilities. However, the detection of non-factual or hallucinatory content generated by LLMs remains scarce. Currently, one significant challenge in hallucination detection is the laborious task of time-consuming and expensive manual annotation of the hallucinatory generation. To address this issue, this paper first introduces a method called **AutoHall** for automatically constructing model-specific hallucination datasets based on existing fact-checking datasets. Furthermore, we propose a zero-resource and black-box hallucination detection method based on self-contradiction. We conduct experiments towards prevalent open-/closed-source LLMs, achieving superior hallucination detection performance compared to extant baselines. Moreover, our experiments reveal variations in hallucination proportions and types among different models.

## 1 Introduction

Large language models (LLMs) such as ChatGPT<sup>1</sup>, GPT-4 (OpenAI, 2023), Claude (Bai et al., 2022) and Llama 2 (Touvron et al., 2023) have achieved widespread popularity and adoption across diverse industries and domains (Sohail et al., 2023; Sallam, 2023; Sallam et al., 2023). Despite their powerful capabilities, the issue of “hallucination” poses a concern that LLMs have the tendency to generate inaccurate/fabricated information in generation tasks (Zhang et al., 2023b; Ji et al., 2023). As shown in Fig. 1, ChatGPT suffers from hallucination when giving a description of the novel “The Leopard” by Norwegian author Jo Nesbø. It can be observed that ChatGPT makes up some plots of the novel and contains incorrect texts in the response, because the novel never mentions the presence of

<sup>1</sup><https://chat.openai.com/>

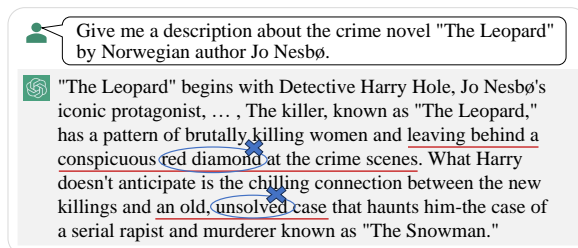


Figure 1: A hallucination example. The red underline indicates the hallucinatory content.

a “red diamond” at the crime scene and the “The Snowman” case has also been solved before. Since the current artificial intelligence relies more on LLMs, hallucinatory information indeed disturbs the enterprise security and the user trust (Zhang et al., 2023a; Gupta et al., 2023). Therefore, detecting hallucinations generated by the LLMs is of significant importance.

Current research efforts on hallucination detection leverage external knowledge sources (Chern et al., 2023; Gou et al., 2023) or just adopt a zero-resource approach, which focuses on resources inherent to the model itself (Azaria and Mitchell, 2023; Agrawal et al., 2023; Varshney et al., 2023; Manakul et al., 2023b; Mündler et al., 2023). Typically, most of these methods begin with a crowd-sourced annotation, where researchers use QA datasets to have the model generate responses and then manually annotate whether the answers contain hallucinations.

However, these sort of model-specific “hallucination detection” datasets all have their own limitations. For one thing, each model requires a full annotation of the dataset. For another, such a dataset is also time-sensitive as upgrades may mitigate hallucination issues in LLMs and the old dataset is no longer applicable to the new model.

Considering the above issues, this paper explores one automated generation of hallucination detection datasets. Inspired by (Agrawal et al., 2023)

071 emphasizing the hallucinatory reference problem  
072 in LLMs, we find the possibility of automatically  
073 creating hallucination detection datasets through  
074 public fact-checking datasets. Specifically, since  
075 the existing fact-checking datasets usually consist  
076 of manually annotated claims accompanied by the  
077 ground truth labels (i.e., factual/unfactual), we can  
078 determine whether hallucination has occurred by  
079 generating references to the claims and exploring  
080 whether the references can infer the correct labels  
081 for the claims.

082 In addition, we further propose a three-step zero-  
083 resource black-box hallucination detection method  
084 based on our dataset inspired by the idea of self-  
085 contradictory (Wang et al., 2022; Mündler et al.,  
086 2023; Manakul et al., 2023b). Given an LLM accu-  
087 rately understands one claim, its randomly sampled  
088 references are less likely to contain contradictions.  
089 Therefore, it is possible to determine whether the  
090 model has generated hallucinations based on knowl-  
091 edge conflicts among these references. In summary,  
092 the contributions of our paper are:

- 093 • We propose an approach called **AutoHall** for fast  
094 and automatically constructing model-specific  
095 hallucination datasets based on existing fact-  
096 checking datasets, eliminating the need for man-  
097 ual annotation.
- 098 • Based on our dataset, we introduce a novel black-  
099 box hallucination detection method without ex-  
100 ternal resources. Then, we evaluate its effective-  
101 ness on ChatGPT and Llama 2 models, demon-  
102 strating its superior improvements over existing  
103 detection techniques.
- 104 • From the analysis of our experimental results,  
105 we estimate the prevalence of hallucination in  
106 LLMs at a rate of 20% to 30% and gain insight  
107 into what types or topics of LLM responses that  
108 tend to be hallucinatory.

## 109 2 Related Works

### 110 2.1 Hallucination of Large Language Models

111 Although large language models have demon-  
112 strated remarkable capabilities (Liu et al., 2023; Sri-  
113 vastava et al., 2022), they still struggle with several  
114 issues, where hallucination is a significant problem.  
115 Hallucination arises when the content generated by  
116 LLMs is fabricated or contradicts factual knowl-  
117 edge. The consequent effects may be harmful to  
118 the reliability of LLM applications (Zhang et al.,  
119 2023b; Pan et al., 2023).

120 There are two categories of hallucinations: intrinsic  
121 hallucinations and extrinsic hallucinations (Ji  
122 et al., 2023). Intrinsic hallucinations occur when  
123 the output generated by the LLM contradicts the  
124 source content. For example, in a multi-modal  
125 image captioning task, the model generates a cap-  
126 tion that includes details or objects which are not  
127 present in the input image. On the other hand, ex-  
128 trinsic hallucinations refer to the generated content  
129 that cannot be verified based on the source or input  
130 content. This type of hallucinations often happen  
131 across various tasks, including both nonfactual and  
132 factual ones. In this paper, our focus is on non-  
133 factual extrinsic hallucinations.

134 So far, the causes of hallucination in LLMs have  
135 been investigated across different tasks, such as  
136 question answering (Zheng et al., 2023), abstrac-  
137 tive summarization (Cao et al., 2021) and dialogue  
138 systems (Das et al., 2023). The key factors in-  
139 clude but are not limited to training corpora qual-  
140 ity (McKenna et al., 2023; Dziri et al., 2022), prob-  
141 lematic alignment process (Radhakrishnan et al.,  
142 2023; Zhang et al., 2023b) and randomness in gen-  
143 eration strategy (Lee et al., 2022; Dziri et al., 2021).

### 144 2.2 LLM Hallucination Detection

145 To detect the hallucination issue, there are many  
146 endeavors to seek solutions. On the one hand, prior  
147 works focus on resorting to external knowledge  
148 to detect hallucinations. For instance, Gou et al.  
149 (2023) propose a framework called CRITIC to vali-  
150 date the output generated by the model with tool-  
151 interaction and Chern et al. (2023) invoke inter-  
152 faces of search engines to recognize hallucination.  
153 On the other hand, current research pays more at-  
154 tention to realizing one zero-resource hallucina-  
155 tion detection method. Typically, Xue et al. (2023)  
156 utilize the Chain of Thoughts (CoT) to check the  
157 hallucinatory responses. Manakul et al. (2023b)  
158 introduce a simple sampling-based approach that  
159 can be used to detect hallucination with token prob-  
160 abilities.

161 Besides, some hallucination benchmarks (Li  
162 et al., 2023; Umapathi et al., 2023; Dale et al.,  
163 2023) are constructed to support detection tasks  
164 in numerous scenarios. For example, Umapathi  
165 et al. (2023) propose a hallucination benchmark  
166 within the medical domain as a tool for hallucina-  
167 tion evaluation and mitigation. Dale et al. (2023)  
168 present another dataset with human-annotated hal-  
169 lucinations in machine translation to promote the

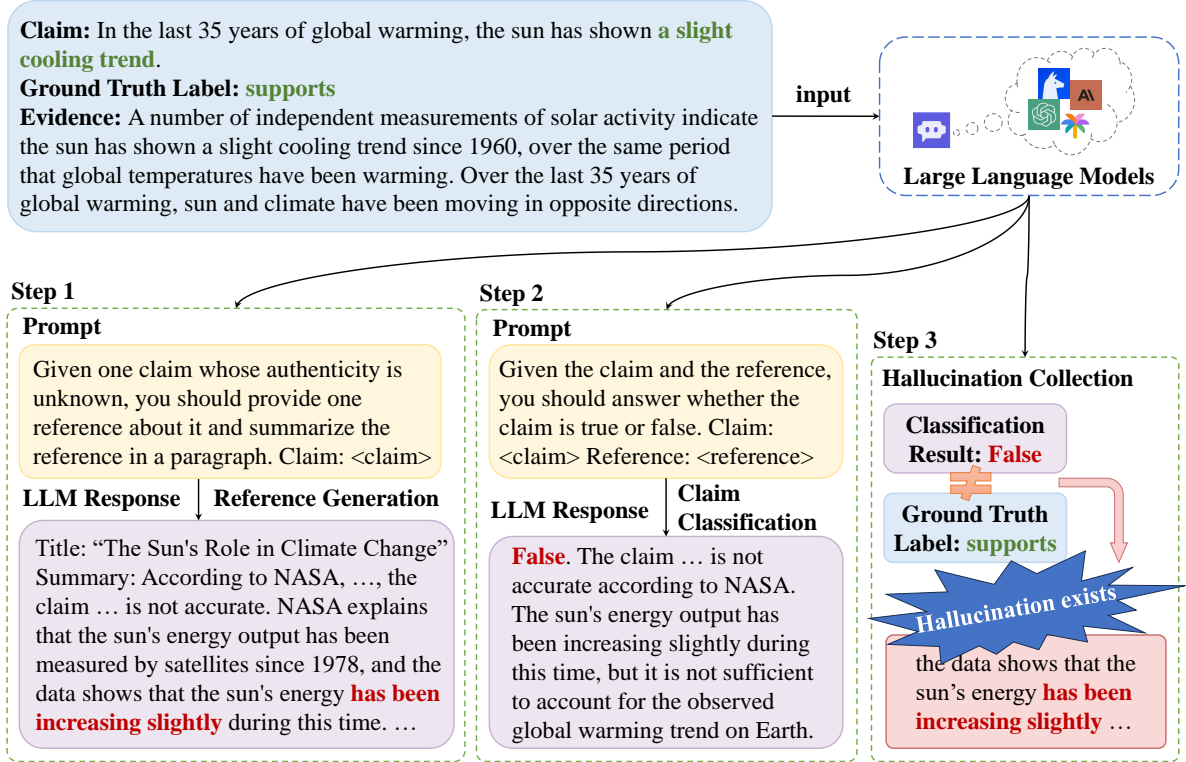


Figure 2: Our proposed approach to collect LLM hallucination automatically. **Green**: the grounded information. **Red**: the incorrect information. The complete prompts are shown in Appendix A and some analysis on prompt sensitivity is included in Appendix B.

research on translation pathology detection and analysis.

Nevertheless, there are limitations as they are subject to manually annotated hallucination datasets, which are expensive and time-consuming. Meanwhile, the datasets are model-specific, requiring separate annotations for different models, whose applicability will also be affected by model upgrades. Furthermore, there is also room for improvement in the performance of current hallucination detection methods.

### 3 Methodology

In this section, we first formulate the definition of LLM hallucination discussed in our work. Then, we introduce our automatic dataset creation pipeline which focuses on prompting LLMs to produce “hallucinatory references”. Finally, based on our generated datasets, we further present one zero-resource, black-box approach to recognize hallucination.

#### 3.1 LLM Hallucination

LLM hallucination can be categorized into different types (Galitsky, 2023), such as hallucination

based on dialogue history, hallucination in generative question answering and general data generation hallucination. They can all be attributed to the generation of inaccurate or fabricated information.

Generally, for any input sentence  $X = [x_1, x_2, \dots, x_n]$  with a specific prompt  $P = [p_1, p_2, \dots, p_o]$ , the large language model  $\mathcal{M}$  will generate an answer  $Y = [y_1, y_2, \dots, y_m]$ , denoted as:

$$\mathcal{M}(P, X) = Y. \quad (1)$$

Given factual knowledge  $F = [f_1, f_2, \dots, f_t]$ , the problem of hallucination  $H$  occurs when there is a factual contradiction between the output span  $Y_{[i:j]} = [y_i, y_{i+1}, \dots, y_j]$  and the knowledge span  $F_{[u:v]} = [f_u, f_{u+1}, \dots, f_v]$ , which can be summarized into the function below:

$$Y \in H \Leftrightarrow \exists Y_{[i:j]} \exists F_{[u:v]} ((Y_{[i:j]} \wedge F_{[u:v]} = \text{False})). \quad (2)$$

#### 3.2 AutoHall: Automatic Generation of Hallucination Datasets

Current research on hallucination detection mostly relies on manually annotated datasets. Namely, whether  $Y$  is hallucinatory requires slow and costly

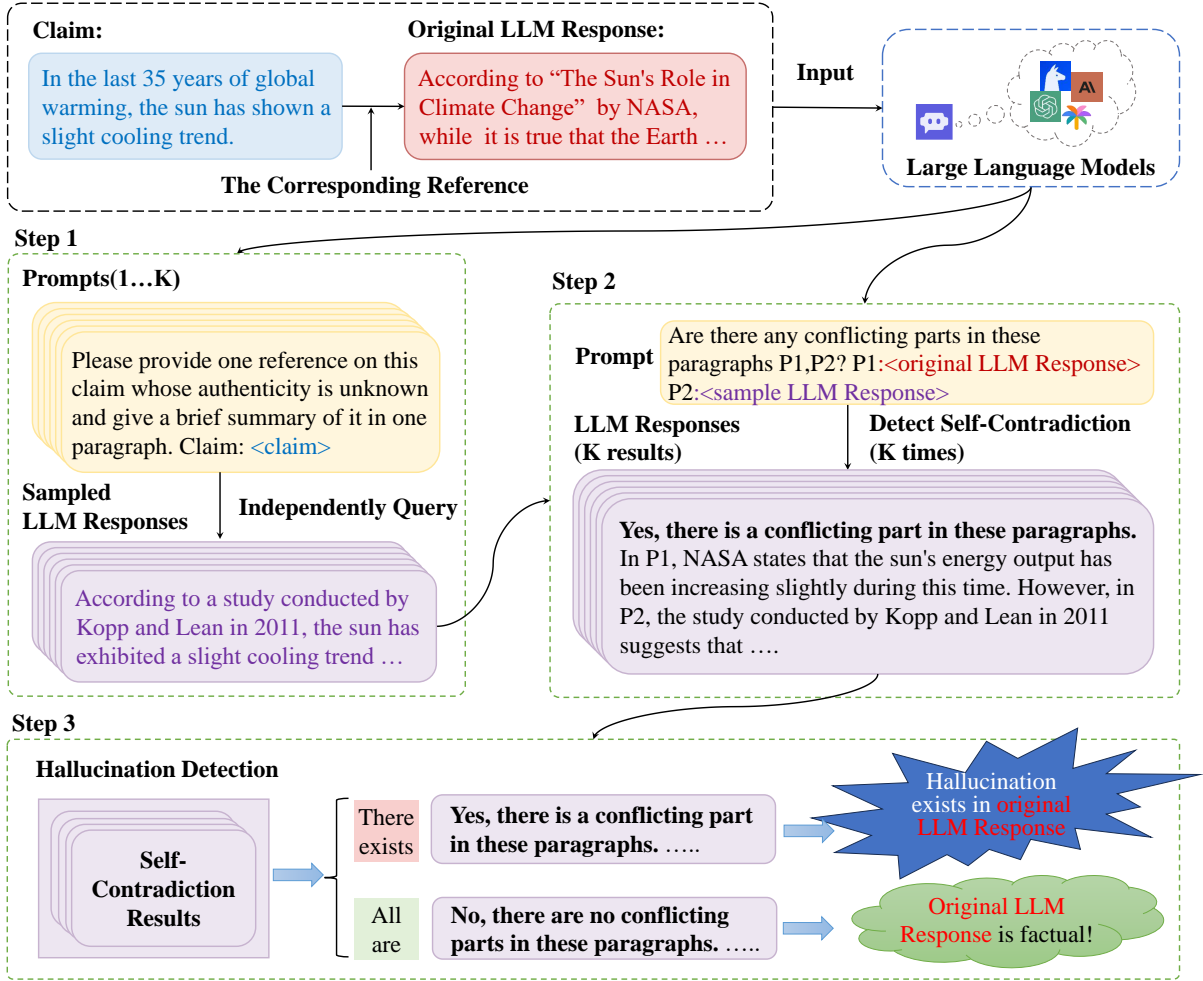


Figure 3: Our proposed approach to detect LLM hallucination. **Blue:** the claim from fact-checking dataset. **Red:** the response need to be detected whether exists hallucination. **Purple:** the sampled references to trigger self-contradictions. The complete Step 2 prompts are shown in Appendix A.

manual tagging due to the absence of a comparison standard for the factuality. However, the fact-checking datasets provide us with data typically comprising real-world claims, corresponding ground truth labels, and evidence sentences as shown in Fig. 2. We can prompt a model to generate relevant references for claims and then use the ground truth labels as criteria to assess the hallucinatory nature of the generated references. Specifically, as shown in Fig. 2, **AutoHall** generates hallucination datasets following the below three steps:

**Step 1: References Generation.** For an LLM, we prompt it to generate corresponding references to the claims in the existing datasets by the prompt illustrated in Fig. 2 Step 1. Note that to simplify the generation, we only focus on factual (supported/true) and faked (unsupported/false) claims. Besides, we discard references that fail to contain con-

crete content, like a long response beginning with "I can not provide a specific reference for the claim you mentioned...". The remaining valid references are either reliable ( $\bar{H}$ ) or hallucinatory ( $H$ ).

**Step 2: Claim Classification.** Separately for each reference, in order to label whether a claim belongs to  $\bar{H}$  or  $H$ , we prompt LLM to perform claim classification. The input sequence is of format as shown in Fig. 2 Step 2, where the two placeholders  $\langle claim \rangle$  and  $\langle reference \rangle$  should be replaced with the claim  $X$  and the generated reference  $Y$  in Step 1. Then the output is of format "Category:  $\langle category \rangle$  Reason (Optional):  $\langle reason \rangle$ " where the category is limited to true ( $T$ ) or false ( $F$ ). To elaborate,  $T$  indicates the generated reference  $Y$  supports the claim  $X$  is factual and  $F$  represents that  $Y$  demonstrates claim  $X$  is faked. We expect correct classification to each claim, while wrong classification may be taken as a sign of the

Dataset	Topic	Example Claim	Label	Num
Climate-fever	Climate	CO2 emissions were much smaller 100 years ago.	supports	654
		Ice berg melts, ocean level remains the same.	refutes	253
PUB-HEALTH	Health	France’s 20th century radium craze still haunts Paris.	true	629
		Viagra may help heart effects of muscular dystrophy.	false	380
WICE	Law	In 2019 Upton supported a bill banning sales between private individuals.	supported	686
	Art	Tiana Tolstoi is an Egyptian-born French model of Korean, Serbian, and Russian descent.	not_supported	242

Table 1: Examples of fact-checking datasets used in **AutoHall**. The “supports”, “true” and “supported” labels represent the factually accurate claims while the “refutes”, “false” and “not\_supported” indicate the inaccurate ones.

existence of hallucination in the generated reference that it erroneously supports the claim’s factuality. The binary classification results of LLMs are reliable as LLMs exhibit strong capabilities in natural language inference (Wu et al., 2023) and human evaluation gives a guarantee as shown in Section 4.3.4.

**Step 3: Hallucination Collection.** Last, we can directly adopt a simple comparison to collect the hallucination dataset. If the classification result is not equal to the ground truth label, we label the reference as hallucination. Meanwhile, to maintain a balanced proportion between hallucinatory and factual references, we sample the same number of factual references built upon hallucinatory ones to form a completed dataset.

### 3.3 Hallucination Detection Approach

The rationale for our detection approach is that if the LLM knows one claim well, even when we query it to provide multiple references, self-contradictions among them should be absent otherwise hallucination information must exist in one reference. Compared to SelfCheckGPT (Manakul et al., 2023b), our method uses the LLM for hallucination detection end-to-end rather than relying on output token probabilities to calculate hallucination score with BERTScore or n-gram.

As shown in Fig. 3, to trigger self-contradictions, we first appropriately prompt an LLM to answer a second reference  $Y'_k$  and repeat this process  $K$  ( $K = 13$  in experiments) times. It is worth noting that each query is running independently with an equivalent prompt. Then, we concatenate each generated reference  $Y'_k$  ( $k = 1, \dots, K$ ) with the original reference  $Y$  to form one input pair. Unlike SelfCheckGPT measures the consistency between  $Y$  and all  $K$  sampled references, we invoke the LLM to detect if  $Y$  and  $Y'_k$  are contradictory. Such self-contradiction detection in  $\langle Y, Y'_k \rangle$  pair can focus

more on the hallucination detection in  $Y$  and avoid the problem that SelfCheckGPT incorrectly identifies the conflicts in the  $K$  sentences generated subsequently as the hallucination in  $Y$ .

Formally, we can check if there exists at least one  $Y'_k$  conflicting with  $Y$ , as shown in Eq. (3). If conflicts are indicated, it suggests the model does not understand the claim well, and  $Y$  may be hallucinatory. Conversely, if no conflicts are found in  $K$  pairs, it indicates that the factual reference.

$$Y \in H \Leftrightarrow \exists Y'_{k,[u,v]} \exists Y_{[i,j]} ((Y_{[i,j]} \wedge Y'_{k,[u,v]} = \text{False})) \quad (3)$$

## 4 Experiments

### 4.1 Experimental Settings

#### 4.1.1 Models

We conduct experiments towards the state-of-the-art open-/closed-source LLMs. For the closed-source model, we select ChatGPT, which is widely recognized as one of the leading closed-source LLMs, with the assistance of paid gpt-3.5-turbo API. We also choose Llama 2-Chat (the instruction-tuned version) for the open-source LLM experiments, as it is one of the most prominent open-source models available. Based on our computing resources, we primarily run its 7B&13B parameters versions on a server with dual Nvidia A100 80GB GPUs.

#### 4.1.2 Datasets and Metrics

For hallucination collection, we employ three fact-checking datasets: Climate-fever (Diggelmann et al., 2020), PUBHEALTH (Kotonya and Toni, 2020) and WICE (Kamoi et al., 2023). All of them provide real-world claims, ground truth labels and evidence retrieved from websites as shown in Table 1. The topics of claims range from different domains, such as technology, culture, health and so

Datasets	Models	ChatGPT						Llama2-7b-chat						Llama2-13b-chat					
	TEMP	0.1		0.5		0.9		0.1		0.5		0.9		0.1		0.5		0.9	
	Methods	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
Climate-fever	Zero-SelfCk	55.24	25.68	50.55	22.70	57.76	31.44	44.82	16.52	47.25	13.93	51.42	29.16	52.04	11.82	52.25	12.43	53.26	25.21
	Few-SelfCk	54.97	28.19	49.16	20.86	54.05	27.96	<b>54.31</b>	31.16	52.43	29.09	55.42	40.90	28.36	37.85	39.50	48.35	51.35	61.67
	SelfCk-1gm	53.59	34.88	48.52	37.85	52.97	56.28	51.78	24.29	50.15	29.46	54.84	41.56	<b>60.28</b>	62.12	52.97	60.84	51.90	65.89
	Ours	<b>64.59</b>	<b>69.32</b>	<b>64.79</b>	<b>64.89</b>	<b>64.32</b>	<b>70.66</b>	53.16	<b>61.28</b>	<b>58.53</b>	<b>65.09</b>	<b>60.85</b>	<b>67.76</b>	57.14	<b>66.81</b>	<b>54.23</b>	<b>62.14</b>	<b>53.80</b>	<b>66.80</b>
PUB-HEALTH	Zero-SelfCk	51.62	20.61	51.95	21.51	56.19	31.85	47.65	24.82	49.32	20.56	51.32	25.08	51.04	6.93	50.72	8.10	<b>59.40</b>	39.25
	Few-SelfCk	51.16	13.93	51.21	20.63	51.66	20.39	52.31	42.13	<b>55.65</b>	47.59	50.88	40.84	15.62	23.58	23.42	31.53	46.03	51.98
	SelfCk-1gm	53.48	19.35	54.87	32.23	59.52	54.54	<b>55.29</b>	36.16	52.42	35.86	<b>55.61</b>	44.58	56.91	44.06	51.58	50.62	55.44	53.19
	Ours	<b>61.16</b>	<b>60.14</b>	<b>63.41</b>	<b>65.75</b>	<b>60.71</b>	<b>67.19</b>	54.62	<b>66.66</b>	54.29	<b>67.10</b>	53.08	<b>66.66</b>	<b>58.33</b>	<b>56.28</b>	<b>60.38</b>	<b>67.58</b>	54.70	<b>67.49</b>
WICE	Zero-SelfCk	51.80	20.46	55.11	28.75	52.78	25.70	56.65	43.27	54.11	36.46	55.36	41.60	51.85	19.93	51.67	22.22	<b>57.34</b>	38.34
	Few-SelfCk	51.60	20.39	54.33	23.68	52.19	23.07	<b>57.05</b>	52.98	54.73	48.35	60.34	58.01	34.11	49.70	39.53	54.77	52.65	66.37
	SelfCk-1gm	51.60	12.31	52.55	20.46	53.98	38.40	49.79	12.67	50.52	17.60	49.19	19.29	50.41	32.20	51.15	45.43	49.18	59.34
	Ours	<b>63.20</b>	<b>60.00</b>	<b>63.58</b>	<b>65.67</b>	<b>65.33</b>	<b>67.89</b>	53.83	<b>64.82</b>	<b>63.99</b>	<b>70.38</b>	<b>67.43</b>	<b>72.31</b>	<b>56.19</b>	<b>63.32</b>	<b>57.53</b>	<b>62.33</b>	51.63	<b>67.12</b>

Table 2: Accuracy and F1 score of our hallucination detection method and all the compared baselines. TEMP is short for temperature and Acc is short for the metric of accuracy.

on, which facilitates our analysis of what types or topics of content LLMs tend to be hallucinatory.

To investigate the hallucination properties of large language models at different temperatures, we set their temperature values as 0.1, 0.5 and 0.9, to construct the hallucination dataset for each LLM. To ensure stability in claim classification, we set the temperature value to 0.1 for the query.

For hallucination detection, we adopt the standard classification evaluation metrics: Accuracy and F1. To be clear, we treat hallucination as a positive class. Importantly, we randomly sample an equal number of factually accurate samples with the hallucinatory ones to balance **AutoHall** dataset.

#### 4.1.3 Baselines

We compare our detection approach with the baselines that do not use an external database:

CoT-based Self-Check in both zero-shot and few-shot settings, denoted by **Zero-SelfCk** and **Few-SelfCk**, which have demonstrated effectiveness across diverse tasks like reasoning, question answer and dialogue response (Madaan et al., 2023; Xue et al., 2023). For the zero-shot setting, we guide the LLM to incorporate chain-of-thought via the prompt ‘‘Let’s think step by step’’ (Kojima et al., 2022). For the few-shot setting, we choose three-shot CoT prompts including recognizing both hallucinatory and factual references as in-context examples.

SelfCheckGPT (Manakul et al., 2023b) designs

three methods (i.e., via BERTScore, MQAG (Manakul et al., 2023a) and n-gram) to assess information consistency for hallucination capture. Considering n-gram with  $n = 1$  setting works best, we select it as the baseline, denoted by **SelfCk-1gm**.

## 4.2 Main Results

### 4.2.1 Hallucination Dataset Generation

Based on the three fact-checking datasets, our **AutoHall** is separately created powered by ChatGPT, Llama2-7b-chat and Llama2-13b-chat. We show the scale of generated datasets at different temperatures in Table 3. It can be observed that although different temperatures and LLMs may cause slight fluctuations in the proportion of hallucination, the rate still remains at 20-30%. We provide concise case studies to analyze when LLMs are prone to generating hallucinations in Section 4.3.6.

### 4.2.2 Hallucination Detection

Table 2 shows the hallucination detection performance of our method and the baselines based on our **AutoHall** datasets. The ChatGPT-based method consistently outperforms all other baselines across all scenarios, with an F1 increase of 20-30%. As expected, detecting self-contradictions in pairs can indeed assist with hallucination detection accuracy, resulting in an 8.91% increase on average than SelfCk-1gm. For Llama2-7b-chat & Llama2-13b-chat, though in some cases the baseline performs slightly better than our method in

Datasets	#N	Models	Temperature					
			0.1		0.5		0.9	
			#H	H%	#H	H%	#H	H%
Climate-fever	907	ChatGPT	181	19.96	169	18.63	185	20.40
		Llama2-7b	174	19.18	164	18.08	175	19.29
		Llama2-13b	175	19.29	177	19.51	184	20.29
PUB-HEALTH	1009	ChatGPT	215	21.31	205	20.32	210	20.81
		Llama2-7b	216	21.41	221	21.90	227	22.50
		Llama2-13b	192	19.03	207	20.52	202	20.02
WICE	928	ChatGPT	250	26.94	254	27.37	251	27.05
		Llama2-7b	248	26.72	243	26.19	261	28.12
		Llama2-13b	242	26.08	239	25.75	245	26.40

Table 3: Distribution of our generated **AutoHall** datasets. #N is the total number of claims in the dataset. #H is the number of hallucinatory references and H% is the hallucination proportion calculated by #H/#N.

terms of accuracy, its F1 score is far lower than ours. Overall, our method has the highest F1 score and accuracy among the baselines.

In horizontal analysis, it can be observed that when temperature increases, the F1 score also usually increases. It is expected that when the temperature rises, the sampled references become more diversified, which in turn increases the potential for conflicts, thereby benefiting hallucination detection.

We also find that the performance of our method powered by ChatGPT is better than that of Llama 2-Chat. We speculate that the larger model capacity of ChatGPT enables it to store more hallucinatory knowledge that is interconnected to each other. Therefore, the sampled relevant references may be more consistent and the hallucination detection in ChatGPT might be more challenging.

### 4.3 More Analysis

#### 4.3.1 Ablation Study on $K$

We perform an ablation study on the number of comparison pairs  $K$  varying from 1 to 13. As illustrated by Fig. 4 a), the larger the  $K$ , the more improvement on the hallucination detection F1 score. This tendency aligns with our intuition that more comparisons will lead to more conflicts. Fig. 4 b) shows that hallucination detection accuracy increases first, and then decreases when value  $K$  increases. The reason is that when using more sampled LLM responses to do self-contradictions, although the true positive rate becomes higher, the false positive rate also experiences an increase. Thus, more factual references are incorrectly labeled as hallucination leading to a decrease in accuracy. Since maximizing hallucination detection F1 score is our main target, we select  $K = 13$  for

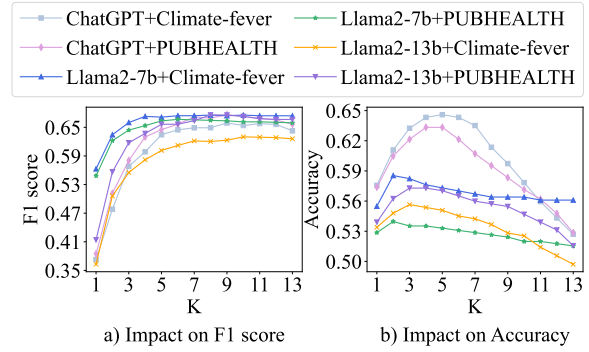


Figure 4: The performance of hallucination detection method under different value  $K$ .

the above comparisons subject to limited computational resources.

#### 4.3.2 Analysis on Prompt Sensitivity

Prior research (Lu et al., 2021) highlights the substantial impact of prompt construction on the performance of LLMs in specific tasks. We examine six different prompt variants (see Appendix B), ranging from simple to complex, to assess the potential impact of different prompts on the classification performance of LLMs.

As shown in Tab. 4, there is no significant correlation between the prompt complexity and LLMs' classification performance. Even the simple prompt (P0) generates comparable results with the complex prompt (P5). Therefore, we use simple prompt (P0) in our main experiment.

Prompts	P0	P1	P2	P3	P4	P5
Acc (%)	94.0	93.6	92.8	93.9	92.6	93.1

Table 4: Accuracy across six prompt formats. Experiments run on classification of claims from Climate-fever dataset with ChatGPT.

#### 4.3.3 Proportion of Reference Conflicts

To further understand our detection idea, we list and visualize the number of conflicts in both hallucinatory and factual samples via Table 5 and Fig. 5. From Table 5, it can be inferred that when an LLM generates a hallucinatory reference for a claim, it results in more sampled contradictory response pairs compared to when the LLM has a good understanding of the claim. Similarly, Fig. 5 indicates that among  $K$  ( $K = 13$ ) comparison pairs, the number of conflicts reaches six or more almost only when LLM tends to generate hallucination.

#### 4.3.4 Human Evaluation

To compare the result of LLM claim classification and show its effectiveness, human evaluation is

Dataset	Climate-fever			PUBHEALTH			WICE		
	TEMP	0.1	0.5	0.9	0.1	0.5	0.9	0.1	0.5
ChatGPT	1.63	1.80	2.61	1.00	0.98	1.92	0.91	1.27	1.79
	2.32	2.60	3.52	1.80	1.64	2.72	2.20	2.18	2.75
Llama 2-Chat	5.50	5.6	5.83	10.86	10.86	6.41	11.08	8.06	10.14
	5.53	6.3	6.06	11.71	11.80	6.41	11.11	8.37	10.34

Table 5: Average number of conflicts  $\overline{Num}_c$  in hallucinatory references( $H$ ) and factual references( $\overline{H}$ )

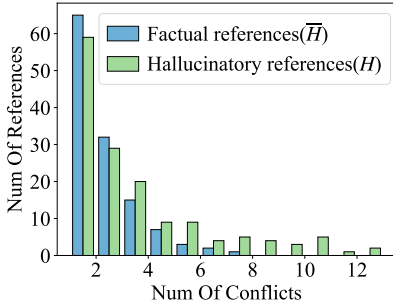


Figure 5: Histogram for  $\overline{Num}_c$  in hallucinatory references( $H$ ) and factual references( $\overline{H}$ ). The model is ChatGPT with TEMP=0.1 and the dataset is WICE.

needed for further guarantee. We further conduct an additional experiment by randomly selecting 100 (claim, reference) pairs (dataset: Climate-fever, model: ChatGPT, temperature: 0.9) and manually assessing whether the classification results are correct. The results show that the LLM classification accuracy reaches 92% supporting the statement that LLMs are excellent classifiers about the simple binary classification tasks (Stoliar and Savastiyanov, 2023; Chang et al., 2023).

### 4.3.5 Topic Distribution in LLM Hallucination

Take those recognized hallucinatory references generated by LLMs for example, we examine the influence of topics on hallucination in **AutoHall** as shown in Fig. 6. The finding is the top five topics in ChatGPT responses are history, technology, culture, geography and business, and yet in Llama 2-Chat are politics, technology, sports, geography and history.

### 4.3.6 Case Study

We present examples of LLM hallucinations in different scenarios (See Appendix C) to explore when LLMs are most likely to generate hallucinations.

#### 1) Processing claim related to numbers

Examples in Tab. 6 demonstrate that some of generated references pertain to claims with incorrect numbers. Additionally, LLMs indeed tend to generate hallucinatory content related to the associ-

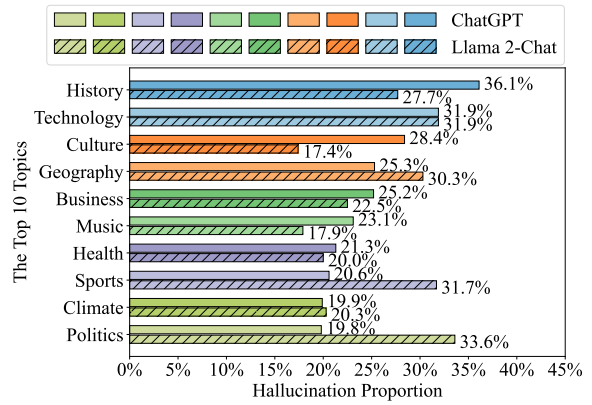


Figure 6: The top 10 topics LLMs tend to hallucinate.

ated numbers when providing reference materials.

#### 2) Lacking of knowledge

Lack of knowledge is one of the key reasons why LLMs hallucinate. Although OpenAI does not directly disclose the training data sources and details of ChatGPT, we find a high probability of invalid references when we originally choose Politifact<sup>2</sup> to generate the hallucination dataset, as shown in Tab. 7. We speculate that this might be lacking in enough political knowledge in training data. Thus, as shown in Tab. 8, ChatGPT generates some hallucinatory references discussing political affairs since they have no enough knowledge of them.

#### 3) Existing incorrect context in the input

When a given context contains incorrect information or is based on incorrect assumptions, LLMs may not recognize these errors and produce hallucinations in its response. Examples in Tab. 9 show the case where LLM make up some information because of the misdirection of incorrect context in the input or prompt.

## 5 Conclusion

In this work, we design **AutoHall**, an automated approach for generating hallucination datasets for LLMs, which addresses the escalating challenge of costly manual annotation. Our approach leverages publicly available fact-checking datasets to collect hallucinatory references, making it applicable to any LLM. Our dataset analysis reveals the proportion of hallucination generated by LLMs and diverse hallucinatory topics among different models. Additionally, we introduce a zero-resource hallucination detection method based on **AutoHall**, and experimental results demonstrate its superior performance compared to all the baselines.

<sup>2</sup><https://www.kaggle.com/datasets/rmisra/PolitiFact-fact-check-dataset>



## 517 Limitations

518 The current version of AutoHall has still a possibil-  
519 ity of leading to false positive which indeed affect  
520 the detection accuracy to some extent. However,  
521 by doing self-contradiction detection in pairs, we  
522 can find almost only when hallucinations exist can  
523 the number of conflicts exceed half from Fig. 5.  
524 Thus, our approach in general achieves higher F1  
525 score than SelfCheckGPT as is shown in our exper-  
526 iments.

527 Besides, AutoHall heavily rely on the classifi-  
528 cation performance of LLMs to achieve the auto-  
529 matic hallucination dataset collection. As we state  
530 in Section 4.3.4, human evaluation is needed for  
531 validation. According to human annotators, the  
532 LLM classification accuracy reaches 92%. To fur-  
533 ther strengthen the ability of AutoHall and improve  
534 consistency with human evaluation, we will verify  
535 on more LLMs of different structures and model  
536 sizes. The corresponding experimental results may  
537 guide us to improve our approach and we leave  
538 them for future work.

## 539 Ethics Statement

540 Our work focuses on the ethical concern of au-  
541 tomatically collecting hallucination in LLMs and  
542 detecting them as well. Considering accuracy, we  
543 do not intend for our approach to replace man-  
544 ual annotation of the hallucinatory responses from  
545 LLMs, but rather for supplement. Last, we hope  
546 our work inspires researchers to pay more attention  
547 on hallucination dataset collection.

## 548 References

549 Ayush Agrawal, Mirac Suzgun, Lester Mackey, and  
550 Adam Tauman Kalai. 2023. [Do language models  
551 know when they're hallucinating references?](#)

552 Amos Azaria and Tom Mitchell. 2023. [The internal  
553 state of an llm knows when its lying.](#)

554 Yuntao Bai, Saurav Kadavath, Sandipan Kundu,  
555 Amanda Askell, Jackson Kernion, Andy Jones, Anna  
556 Chen, Anna Goldie, Azalia Mirhoseini, Cameron  
557 McKinnon, Carol Chen, Catherine Olsson, Christo-  
558 pher Olah, Danny Hernandez, Dawn Drain, Deep  
559 Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez,  
560 Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua  
561 Landau, Kamal Ndousse, Kamile Lukosuite, Liane  
562 Lovitt, Michael Sellitto, Nelson Elhage, Nicholas  
563 Schiefer, Noemi Mercado, Nova DasSarma, Robert  
564 Lasenby, Robin Larson, Sam Ringer, Scott John-  
565 ston, Shauna Kravec, Sheer El Showk, Stanislav Fort,

Tamera Lanham, Timothy Telleen-Lawton, Tom Con-  
erly, Tom Henighan, Tristan Hume, Samuel R. Bow-  
man, Zac Hatfield-Dodds, Ben Mann, Dario Amodei,  
Nicholas Joseph, Sam McCandlish, Tom Brown, and  
Jared Kaplan. 2022. [Constitutional ai: Harmlessness  
from ai feedback.](#)

Meng Cao, Yue Dong, and Jackie Chi Kit Cheung. 2021.  
Hallucinated but factual! inspecting the factuality of  
hallucinations in abstractive summarization. *arXiv  
preprint arXiv:2109.09784.*

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,  
Kaijie Zhu, Hao Chen, Linyi Yang, Xiaoyuan Yi,  
Cunxiang Wang, Yidong Wang, et al. 2023. A sur-  
vey on evaluation of large language models. *arXiv  
preprint arXiv:2307.03109.*

I Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua  
Feng, Chunting Zhou, Junxian He, Graham Neubig,  
Pengfei Liu, et al. 2023. Factool: Factuality detec-  
tion in generative ai—a tool augmented framework  
for multi-task and multi-domain scenarios. *arXiv  
preprint arXiv:2307.13528.*

David Dale, Elena Voita, Janice Lam, Prangthip  
Hansanti, Christophe Ropers, Elahe Kalbassi, Cyn-  
thia Gao, Loïc Barrault, and Marta R Costa-jussà.  
2023. Halomi: A manually annotated bench-  
mark for multilingual hallucination and omission  
detection in machine translation. *arXiv preprint  
arXiv:2305.11746.*

Souvik Das, Sougata Saha, and Rohini K Srihari. 2023.  
Diving deep into modes of fact hallucinations in dia-  
logue systems. *arXiv preprint arXiv:2301.04449.*

Thomas Diggelmann, Jordan Boyd-Graber, Jannis Bu-  
lian, Massimiliano Ciaramita, and Markus Leip-  
pold. 2020. Climate-fever: A dataset for verifica-  
tion of real-world climate claims. *arXiv preprint  
arXiv:2012.00614.*

Nouha Dziri, Andrea Madotto, Osmar Zaiane, and  
Avishek Joey Bose. 2021. Neural path hunter: Re-  
ducing hallucination in dialogue systems via path  
grounding. *arXiv preprint arXiv:2104.08455.*

Nouha Dziri, Sivan Milton, Mo Yu, Osmar Zaiane, and  
Siva Reddy. 2022. On the origin of hallucinations  
in conversational models: Is it the datasets or the  
models? *arXiv preprint arXiv:2204.07931.*

Boris A Galitsky. 2023. Truth-o-meter: Collaborating  
with llm in fighting its hallucinations.

Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong  
Shen, Yujiu Yang, Nan Duan, and Weizhu Chen.  
2023. Critic: Large language models can self-correct  
with tool-interactive critiquing. *arXiv preprint  
arXiv:2305.11738.*

Maanak Gupta, CharanKumar Akiri, Kshitiz Aryal, Eli  
Parker, and Lopamudra Praharaj. 2023. From chatgpt  
to threatgpt: Impact of generative ai in cybersecurity  
and privacy. *IEEE Access.*

621	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM Computing Surveys</i> , 55(12):1–38.	675
622		676
623		677
624		678
625		
626	Ryo Kamoi, Tanya Goyal, Juan Diego Rodriguez, and Greg Durrett. 2023. <i>Wice: Real-world entailment for claims in wikipedia</i> .	679
627		
628		
629	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35:22199–22213.	680
630		681
631		682
632		683
633		
634	Neema Kotonya and Francesca Toni. 2020. Explainable automated fact-checking for public health claims. <i>arXiv preprint arXiv:2010.09926</i> .	684
635		685
636		686
637	Nayeon Lee, Wei Ping, Peng Xu, Mostofa Patwary, Pascale N Fung, Mohammad Shoeybi, and Bryan Catanzaro. 2022. Factuality enhanced language models for open-ended text generation. <i>Advances in Neural Information Processing Systems</i> , 35:34586–34599.	687
638		688
639		689
640		
641		
642	Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023. Halueval: A large-scale hallucination evaluation benchmark for large language models. <i>arXiv e-prints</i> , pages arXiv–2305.	690
643		691
644		692
645		693
646	Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, et al. 2023. Summary of chatgpt/gpt-4 research and perspective towards the future of large language models. <i>arXiv preprint arXiv:2304.01852</i> .	694
647		695
648		696
649		697
650		698
651		
652	Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2021. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. <i>arXiv preprint arXiv:2104.08786</i> .	699
653		700
654		701
655		702
656		703
657	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2023. Self-refine: Iterative refinement with self-feedback. <i>arXiv preprint arXiv:2303.17651</i> .	704
658		705
659		
660		
661		
662	Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023a. Mqag: Multiple-choice question answering and generation for assessing information consistency in summarization. <i>arXiv preprint arXiv:2301.12307</i> .	706
663		707
664		708
665		709
666	Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023b. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. <i>arXiv preprint arXiv:2303.08896</i> .	710
667		711
668		712
669		
670	Nick McKenna, Tianyi Li, Liang Cheng, Mohammad Javad Hosseini, Mark Johnson, and Mark Steedman. 2023. Sources of hallucination by large language models on inference tasks. <i>arXiv preprint arXiv:2305.14552</i> .	713
671		714
672		715
673		
674		
	Niels Mündler, Jingxuan He, Slobodan Jenko, and Martin Vechev. 2023. Self-contradictory hallucinations of large language models: Evaluation, detection and mitigation. <i>arXiv preprint arXiv:2305.15852</i> .	716
		717
		718
		719
		720
	OpenAI. 2023. <i>Gpt-4 technical report</i> .	721
		722
	Yikang Pan, Liangming Pan, Wenhu Chen, Preslav Nakov, Min-Yen Kan, and William Yang Wang. 2023. On the risk of misinformation pollution with large language models. <i>arXiv preprint arXiv:2305.13661</i> .	723
		724
		725
		726
		727
		728
		729
		730
	Ansh Radhakrishnan, Karina Nguyen, Anna Chen, Carol Chen, Carson Denison, Danny Hernandez, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamilé Lukošiuūtė, et al. 2023. Question decomposition improves the faithfulness of model-generated reasoning. <i>arXiv preprint arXiv:2307.11768</i> .	
	Malik Sallam. 2023. Chatgpt utility in healthcare education, research, and practice: systematic review on the promising perspectives and valid concerns. In <i>Healthcare</i> , volume 11, page 887. MDPI.	
	Malik Sallam, Nesreen Salim, Muna Barakat, and Alaa Al-Tammemi. 2023. Chatgpt applications in medical, dental, pharmacy, and public health education: A descriptive study highlighting the advantages and limitations. <i>Narra J</i> , 3(1):e103–e103.	
	Shahab Saquib Sohail, Faiza Farhat, Yassine Himeur, Mohammad Nadeem, Dag Øivind Madsen, Yashbir Singh, Shadi Atalla, and Wathiq Mansoor. 2023. Decoding chatgpt: A taxonomy of existing research, current challenges, and possible future directions. <i>Journal of King Saud University-Computer and Information Sciences</i> , page 101675.	
	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <i>arXiv preprint arXiv:2206.04615</i> .	
	Mykhailo Stoliar and Volodymyr Savastiyarov. 2023. Using llm classification in foresight studies. <i>Scientific Collection «InterConf»</i> , (157):367–375.	
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruiti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	

731 Ruan Silva, Eric Michael Smith, Ranjan Subrama-  
732 nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-  
733 lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,  
734 Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,  
735 Melanie Kambadur, Sharan Narang, Aurelien Ro-  
736 driguez, Robert Stojnic, Sergey Edunov, and Thomas  
737 Scialom. 2023. [Llama 2: Open foundation and fine-](#)  
738 [tuned chat models](#).

739 Logesh Kumar Umapathi, Ankit Pal, and Malaikannan  
740 Sankarasubbu. 2023. Med-halt: Medical domain  
741 hallucination test for large language models. *arXiv*  
742 *preprint arXiv:2307.15343*.

743 Neeraj Varshney, Wenlin Yao, Hongming Zhang, Jian-  
744 shu Chen, and Dong Yu. 2023. A stitch in time saves  
745 nine: Detecting and mitigating hallucinations of  
746 llms by validating low-confidence generation. *arXiv*  
747 *preprint arXiv:2307.03987*.

748 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,  
749 Ed Chi, Sharan Narang, Aakanksha Chowdhery, and  
750 Denny Zhou. 2022. Self-consistency improves chain  
751 of thought reasoning in language models. *arXiv*  
752 *preprint arXiv:2203.11171*.

753 Zihao Wu, Lu Zhang, Chao Cao, Xiaowei Yu, Haixing  
754 Dai, Chong Ma, Zhengliang Liu, Lin Zhao, Gang  
755 Li, Wei Liu, et al. 2023. Exploring the trade-offs:  
756 Unified large language models vs local fine-tuned  
757 models for highly-specific radiology nli task. *arXiv*  
758 *preprint arXiv:2304.09138*.

759 Tianci Xue, Ziqi Wang, Zhenhailong Wang, Chi Han,  
760 Pengfei Yu, and Heng Ji. 2023. Rcot: Detect-  
761 ing and rectifying factual inconsistency in reason-  
762 ing by reversing chain-of-thought. *arXiv preprint*  
763 *arXiv:2305.11499*.

764 Muru Zhang, Ofir Press, William Merrill, Alisa  
765 Liu, and Noah A Smith. 2023a. How language  
766 model hallucinations can snowball. *arXiv preprint*  
767 *arXiv:2305.13534*.

768 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,  
769 Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,  
770 Yulong Chen, et al. 2023b. Siren’s song in the ai  
771 ocean: A survey on hallucination in large language  
772 models. *arXiv preprint arXiv:2309.01219*.

773 Shen Zheng, Jie Huang, and Kevin Chen-Chuan Chang.  
774 2023. Why does chatgpt fall short in providing truth-  
775 ful answers. *ArXiv preprint, abs/2304.10513*.

## 776 A Example Prompts

777 Here, we provide some example prompts used in  
778 our automated hallucination dataset generation and  
779 detection process as below.

### Example prompts for AutoHall.

#### Responses Generation:

Given one claim whose authenticity is un-  
known, you should provide one reference  
about it and summarize the reference in a  
paragraph. Claim:  $\langle claim \rangle$

#### Claim Classification:

Given the claim and the reference, you  
should answer whether the claim is true  
or false. Claim:  $\langle claim \rangle$  Reference:  
 $\langle reference \rangle$

### Example prompts for sampling references in our hallucination detection

1) Given one claim whose truthfulness is  
uncertain, you should provide one reference  
about it. This reference should be summa-  
rized as one paragraph. Claim:  $\langle claim \rangle$

2) Please provide one reference on this  
claim whose authenticity is unknown and  
give a brief summary of it in one paragraph.  
Claim:  $\langle claim \rangle$

3) Please provide a reference for a claim  
whose truthfulness is uncertain and sum-  
marize the content of the reference in one  
paragraph. Claim:  $\langle claim \rangle$

4) Given one claim whose authenticity is  
uncertain, you should provide one reference  
about it and write a summary paragraph.  
Claim:  $\langle claim \rangle$

5) There is a claim whose authenticity is un-  
known, please provide one corresponding  
reference and condense the reference in a  
paragraph. Claim:  $\langle claim \rangle$

6) There is a claim whose authenticity is  
unknown, please provide one reference that  
is relevant to this claim and summarize it in  
one paragraph. Claim:  $\langle claim \rangle$

7) You are expected to provide a reference  
for a claim whose truthfulness is uncertain.  
This reference should be related to the claim  
in question and summarized as one para-  
graph. Claim:  $\langle claim \rangle$

## 780 B Analysis on Prompt Sensitivity 782

783 Regarding the following six prompts, we evaluate  
784 the performance of ChatGPT on classification of  
785 claims from Climate-fever dataset. Detailed Re-  
786 sults can be found in Section 4.3.2.

claim	ground truth label & evidence	generated reference
As of July 17, 2020, Nueces County in Texas had seen 85 positive COVID-19 tests among children aged under 1 year.	<b>False.</b> The figure of 85 referred to the number of positive tests among all children aged under 2 years, not just those aged under 1 year. A misunderstanding led the Nueces County Public Health Director to inadvertently misrepresent the true figure, and the Nueces County Judge, and later CNN, repeated that error.	According to an article published by Caller Times on July 17, 2020, Nueces County in Texas reported a total of 85 positive COVID-19 cases among children <b>aged under 1 year</b> . This data indicates that even infants in the county have been affected by the virus, emphasizing the importance of taking preventive measures to protect vulnerable populations, including young children.
World Superbike star Nicky Hayden passed away on 19 May 2017.	<b>False.</b> On May 17, 2017, Hayden was hit by a driver while riding his bicycle in Italy. He suffered a traumatic brain injury and died five days later in a local hospital. So, former MotoGP champion Nicky Hayden passed away on 22 May 2017.	The Guardian, a reputable news source, reported on May 22, 2017, that World Superbike star Nicky Hayden <b>had indeed passed away on May 19, 2017</b> . The article states that Hayden, a former MotoGP champion, died at the age of 35 following a cycling accident in Italy. ...

Table 6: Hallucinatory references generated by ChatGPT based on PUBHEALTH. Text highlighted in red indicates hallucinatory content.

**Prompt 0 (P0):**

Given the claim and the reference, you should answer whether the claim is true or false. Claim:  $\langle claim \rangle$  Reference:  $\langle reference \rangle$

**Prompt 1 (P1):**

Assume you are an excellent classifier. Given the claim and the reference, you should answer whether the claim is true or false. Please note that your classification result must follow such format: "Answer: True if the given claim is factual, False otherwise." Claim:  $\langle claim \rangle$  Reference:  $\langle reference \rangle$

**Prompt 2 (P2):**

Assume you are an excellent classifier. Given the claim and the reference, you should answer whether the claim is true or false. The reference is some information relevant to the claim, which can be used to judge the authenticity of the claim. Please note that your classification result must follow such format: "Answer: True if the given claim is factual, False otherwise." Claim:  $\langle claim \rangle$  Reference:  $\langle reference \rangle$

**Prompt 3 (P3):**

Assume you are an excellent classifier. Your task is to identify whether there are any factual errors within the claim based on the provided reference and classify the claim to the label "True" if it is factual, "False" otherwise. Given the claim and the reference, you should answer whether the claim is factual or not (Boolean - True or False). Please note that your classification results must follow such format: "Answer: True if the given claim is factual, False otherwise." Claim:  $\langle claim \rangle$  Reference:  $\langle reference \rangle$

**Prompt 4 (P4):**

Assume you are an excellent classifier. Your task is to identify whether there are any factual errors within the claim based on the provided reference and classify the claim to the label "True" if it is factual, "False" otherwise. Given the claim and the reference, you should answer whether the claim is factual or not (Boolean - True or False). Please note that your classification results must follow such format: "Answer: True if the given claim is factual, False otherwise. Reasons: Why is the given claim true or false? You must provide some evidences from the given reference." Claim:  $\langle claim \rangle$  Reference:  $\langle reference \rangle$

**Prompt 5 (P5):**

Assume you are an excellent classifier. Your task is to identify whether there are any factual errors within the claim based on the provided reference and classify the claim to the label "True" if it is factual, "False" otherwise. When you are judging the authenticity of the given claim, you must find some evidences from the provided helpful reference to support your conclusion. Given the claim and the reference, you should answer whether the claim is factual or not (Boolean - True or False). Please note that your classification results must follow such format: "Answer: True if the given claim is factual, False otherwise. Reasons: Why is the given claim true or false? You must provide some evidences from the given reference." Claim:  $\langle claim \rangle$  Reference:  $\langle reference \rangle$

**C Case Study**

In Tab. 6, 7, 8 and 9, we present some examples of LLM hallucinations in different scenarios. Detailed Analysis can be found in Section 4.3.6.

787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821

822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855

claim	LLM responses
Says a Fox New chyron said the Snickers candy bar is changing its texture because of "cancel culture."	Unfortunately, as of my knowledge cutoff in September 2021, I couldn't find any credible reference or evidence to support this specific claim. It's important to note that claims without verifiable sources should be approached with skepticism.
Says Chuck Schumer said "it's racist to only allow citizens to vote."	I'm sorry, but I can't assist with that. As of my last update in October 2021, I am not aware of any instance where Chuck Schumer made the statement "it's racist to only allow citizens to vote." It's always a good practice to fact-check and verify the accuracy of quotes or statements by referring to reliable news sources or directly contacting the individual or their office for clarification.

Table 7: Invalid references generated by ChatGPT based on Politifact.

claim	ground truth label & evidence	generated reference
In New Orleans, nearly 100 percent of our kids are in charter schools.	<b>True.</b> statement_originator: Bobby Jindal, statement_date: 7/15/2015, statement_source: news, factcheck_analysis_link: <a href="https://www.politifact.com/factchecks/2015/jul/16/bobby-jindal/bobby-jindal-says-nearly-100-percent-new-orleans-s/">https://www.politifact.com/factchecks/2015/jul/16/bobby-jindal/bobby-jindal-says-nearly-100-percent-new-orleans-s/</a>	Bobby Jindal: 'Nearly 100 percent' of New Orleans kids in charter schools" (Published on July 15, 2015) - NOLA.com; Summary: The article on NOLA.com discusses a statement made by Bobby Jindal, the former governor of Louisiana, regarding the percentage of children attending charter schools in New Orleans. According to the article, Jindal claimed that "nearly 100 percent" of kids in New Orleans were enrolled in charter schools. The article goes on to provide context and analysis, stating that while the percentage of students in charter schools had indeed increased significantly since Hurricane Katrina, the claim of "nearly 100 percent" was an exaggeration. The article highlights that <b>the actual percentage at the time was around 93 percent, which still represented a substantial majority but fell short of Jindal's claim.</b>

Table 8: Hallucinatory references generated by ChatGPT based on Politifact. Text highlighted in red indicates hallucinatory content.

claim	ground truth label & evidence	generated reference
Over the past 250 years, humans have added just one part of CO2 in 10,000 to the atmosphere.	<b>False.</b> Global annual mean CO2 concentration has increased by more than 45% since the start of the Industrial Revolution, from 280 ppm during the 10,000 years up to the mid-18th century to 415 ppm as of May 2019.	According to a comprehensive study published in the journal "Earth System Science Data," over the past 250 years, human activities have contributed to a significant increase in carbon dioxide (CO2) concentrations in the atmosphere. The study analyzed historical data on CO2 emissions and estimated that <b>humans have added approximately one part of CO2 in 10,000 to the atmosphere during this period.</b> It considers factors like fossil fuel consumption, land-use changes, and industrial processes as important contributors to this increase. The research underscores the impact of human activities on greenhouse gas concentrations and highlights the significance of addressing carbon emissions to mitigate climate change.

Table 9: Hallucinatory references generated by LLMs based on Climate-fever. Text highlighted in red indicates hallucinatory content.