

# Investigating and Mitigating the Multimodal Hallucination Snowballing in Large Vision-Language Models

Weihong Zhong<sup>1</sup> Xiaocheng Feng<sup>1,2\*</sup> Liang Zhao<sup>1</sup> Qiming Li<sup>1</sup> Lei Huang<sup>1</sup>  
Yuxuan Gu<sup>1</sup> Weitao Ma<sup>1</sup> Yuan Xu<sup>1</sup> Bing Qin<sup>1,2</sup>

<sup>1</sup>Harbin Institute of Technology

<sup>2</sup>Peng Cheng Laboratory

{whzhong, xcfeng, lzhao, qmli, lhuang, yxgu, wtma, yuanxu, qinb}@ir.hit.edu.cn

## Abstract

Though advanced in understanding visual information with human languages, Large Vision-Language Models (LVLMs) still suffer from multimodal hallucinations. A natural concern is that during multimodal interaction, the generated hallucinations could influence the LVLMs' subsequent generation. Thus, we raise a question: *When presented with a query relevant to the previously generated hallucination, will LVLMs be misled and respond incorrectly, even though the ground visual information exists?* To answer this, we propose a framework called *MMHalSnowball* to evaluate LVLMs' behaviors when encountering generated hallucinations, where LVLMs are required to answer specific visual questions within a curated hallucinatory conversation. Crucially, our experiment shows that the performance of open-source LVLMs drops by at least 31%, indicating that LVLMs are prone to accept the generated hallucinations and make false claims that they would not have supported without distractions. We term this phenomenon *Multimodal Hallucination Snowballing*. To mitigate this, we further propose a training-free method called *Residual Visual Decoding*, where we revise the output distribution of LVLMs with the one derived from the residual visual input, providing models with direct access to the visual information. Experiments show that our method can mitigate more than 24% of the snowballed multimodal hallucination while maintaining capabilities.<sup>1</sup>

## 1 Introduction

Large Vision-Language Models (LVLMs) have shown remarkable abilities in observing and understanding the real world in human languages (Achiam et al., 2023; Zhu et al., 2023; Liu et al., 2023d; Ye et al., 2023b; Dai et al.). However, multimodal hallucinations, in which LVLMs provide

\*Corresponding Author

<sup>1</sup>Resources will be available at <https://github.com/whongzhong/MMHalSnowball>

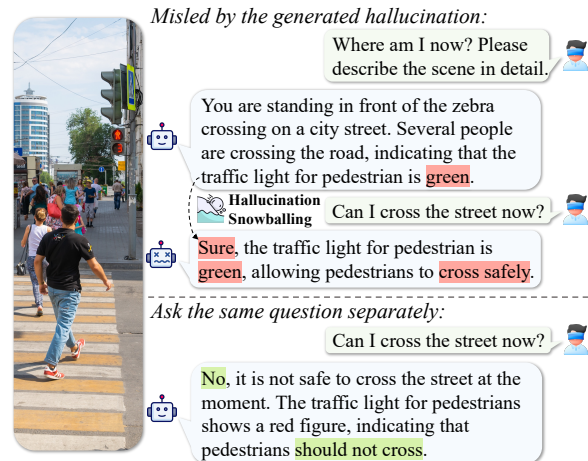


Figure 1: An example of the LVLM assisting a visually impaired person to cross the street. The model is misled by the generated hallucination and mistakenly suggests the user to cross the street, although it can give correct advice independently. Green and red colors highlight the correct answer and hallucinations, respectively.

responses misaligned with the corresponding visual information, remain to be the Achilles' heel (Cui et al., 2023; Kamath et al., 2023; Li et al., 2023b; Liu et al., 2023a; Lu et al., 2023; Rawte et al., 2023; West et al., 2023; Huang et al., 2023).

Previous research has revealed that hallucinations generated by large language models may accumulate due to models' over-commitment to early mistakes, leading to more mistakes that they otherwise would not make (Zhang et al., 2023a; Azaria and Mitchell, 2023; Kang et al., 2023), especially for the user-model interaction scenarios such as conversation (Huang et al., 2022; Tian et al., 2024; Gong et al., 2023). However, the extent to which accumulated multimodal hallucinations mislead LVLMs into generating false claims requires further exploration. In this work, we conducted an investigation into this issue for the first time. As shown in Figure 1, we seek the answer to the question: *When presented with a query relevant to the*

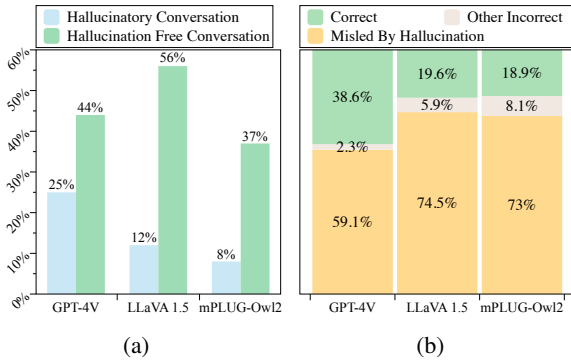


Figure 2: Preliminary explorations on the hallucinations generated by LVLMs given conversational contexts. (a) Response accuracy with or without hallucinatory conversation. (b) Response distribution when asking the question within a hallucinatory conversation. We select question samples that the LVLm can correctly answer without distractions.

previously generated hallucination that contradicts the visual information, can models make the correct judgment when they could have given a correct answer independently? We conduct a preliminary study on GPT-4V (Achiam et al., 2023), LLaVA 1.5 (Liu et al., 2023c), and mPLUG-Owl2 (Ye et al., 2023b). Similar to the setting of Figure 1, given an image, we start a conversation by asking the model to describe the image in detail. When observing hallucinations in the LVLm’s responses, we continue to ask a relevant question according to the model-generated hallucination. In addition, we ask the same question separately to see if the model can answer it correctly without distractions. As demonstrated in Figure 2(a), we find that when the text context contains relevant hallucination, the model performance declines significantly, compared to the model response when asking the same question separately. We further select those question samples that the LVLm can correctly answer separately, and manually identify the response change when asking the same question with the related model-generated hallucinatory context. As Figure 2(b) depicts, we find that more than 59% of the answers are semantically the same as the generated hallucination, indicating that they were misled by the previously generated hallucinations.

To systematically investigate this phenomenon, we propose to identify whether the LVLm is misled by hallucinations via checking if a specific claim is flipped due to previous hallucinations. We design a framework called *MMHalSnowball* to construct hallucinatory visual conversations, where models

are required to answer the question based on the image and the hallucinatory conversation. The result shows that LVLms’ multimodal hallucinations are easy to mislead the later generation because their strong language capabilities make them prone to be over-confident in the hallucinated context, thereby generating false claims that they normally would not support, which we term as *Multimodal Hallucination Snowballing*.

In addition to mitigating this issue, we further proposed a training-free decoding method called *Residual Visual Decoding* (RVD). By residual connecting the visual information and the current user instruction, distributions that emphasizing the visual information are derived to revise the original output distribution. Our RVD achieves more than 24% of improvements in reducing the multimodal hallucination snowballing while maintaining the contextual modeling ability.

## 2 Evaluating the Multimodal Hallucination Snowball Phenomenon

In this section, we design a question-answer task in the conversation scenario, where a model is first asked to describe a picture in detail and then answers a visual question. As shown in Figure 3, we propose the *MMHalSnowball* framework to carefully simulate hallucinatory conversations and evaluate whether the model generates a wrong answer due to the hallucinatory context. Next, we will describe our evaluation framework in detail, including conversation creation, experimental settings, and evaluation metrics. We experimentally analyze the multimodal hallucinations snowball in §2.7. The prompts used are listed in Appendix A.2.

### 2.1 Dataset Source

We use the validation set of the GQA dataset (Hudson and Manning, 2019) as our data source, which contains a balanced aspect of visual questions that focuses on objective perceptual questions. We adopt images, question-answer pairs, and regional description annotations from the Visual Genome (Krishna et al., 2017). Note that we use its balanced validation set to minimize the impact of dataset contamination and language prior.

### 2.2 Hallucination Allocation

To be more practical, we construct hallucinations based on the common types generated by LVLms. Inspired by Wang et al. (2023a); Zhai et al. (2023), we categorize the hallucinations as follows:

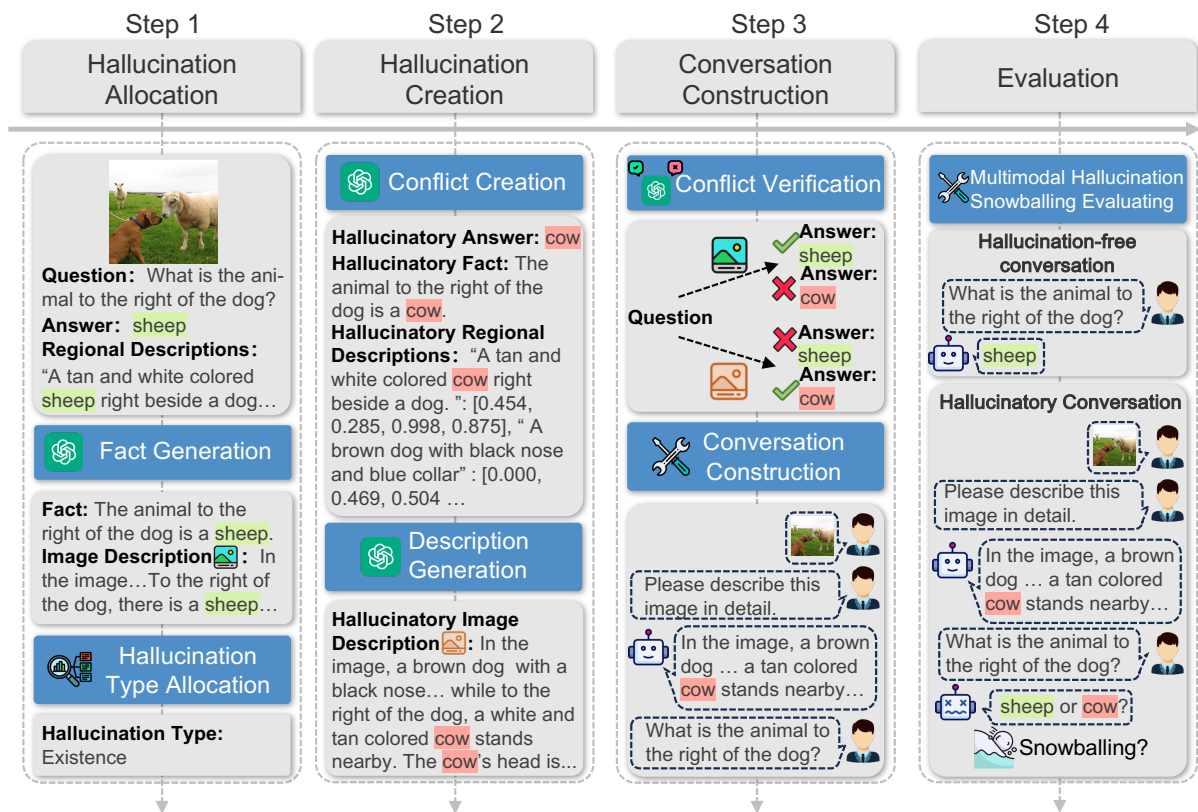


Figure 3: An overview of our *MMHalSnowball* framework for simulating hallucinatory conversations and evaluating LVLMs' behavior in such conversations. In step 1, start with a question-answer pair, we generate a fact, an image description and allocate a proper hallucination type according to the corresponding question-answer pair. In step 2, we utilize the ChatGPT to rewrite a hallucinatory answer based on the allocated hallucination type. We then modify other annotations and generate the corresponding hallucinatory description using ChatGPT. In step 3, after ensuring the hallucinatory answer and descriptions contradict the image content, we construct a conversation that contains the specific hallucination. In step 4, we evaluate the LVLMs' performance gap in two conversation settings to see whether they suffer from multimodal hallucination snowballing. Green and red color highlight the correct answer and hallucinations curated out of it, respectively.

- *Existence Hallucination*, which refers to the incorrect recognition of visible objects in the image or the belief that specific visible objects are absent in the image.
- *Attribute Hallucination*, which refers to the inaccurate characterization of objects and misrepresentations of attributes such as color, shape, size, and actions.
- *Relation Hallucination*, which refers to the inaccurate depiction of the relationships or interactions among objects, including erroneous interaction states, relative positions, and spatial positions of objects relative to the image.
- *Imagination Hallucination*, which refers to the erroneous imagination of objects in the picture that do not appear.

To incorporate hallucinations, we first utilize ChatGPT (OpenAI, 2022) to rewrite a fact sentence

that best describes the question-answer pair. In addition, the annotated regional descriptions and the fact sentence are used to generate an image description. The ChatGPT is prompted to ensure the image description semantically entails the fact sentence. Hallucination can be created by properly modifying the fact sentence. Our goal is to make the answer to the original question no longer correct according to the modified fact sentence. However, not all types of hallucination will make the original answer invalid (e.g. modify the fact sentence "The color of the trousers is blue" to "the color of the bike is blue" introduces an imagination hallucination, but won't invalidate the answer to the question: "What color are the trousers that this boy is wearing in the image?"). To match the hallucination errors in the curated contexts with the corresponding question-answer pairs, We then

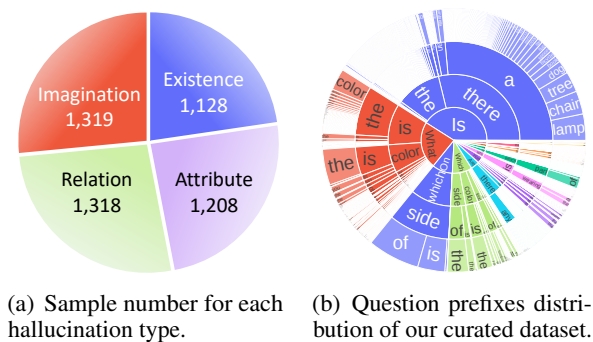


Figure 4: Statistics of our curated dataset.

allocate a proper hallucination type from the above definition to each fact sentence. Appendix A.1 shows details about the rules of allocating proper hallucination types.

### 2.3 Hallucination Creation

In this part, we describe how we utilize the question-answer pair, the fact sentence, and the regional descriptions to generate hallucinatory image descriptions. Rather than directly modifying the fact sentence according to the hallucination type to create hallucinations, we find it more stable to ask the ChatGPT to rewrite a hallucinatory answer that contradicts the original answer. Then, the fact sentence, as well as all the regional descriptions are heuristically modified to hallucinatory ones according to the hallucinatory answer. With the hallucinatory fact sentence and hallucinatory regional descriptions as inputs, the ChatGPT is asked to generate a detailed image description that entails the hallucinatory fact. the original answers  $Y^+ = \{y_1^+, y_1^+, \dots, y_n^+\}$  and the rewritten hallucinatory answers  $Y^- = \{y_1^-, y_1^-, \dots, y_n^-\}$  are kept for evaluation, where  $n$  represents the dataset size.

### 2.4 Conversation Construction

Before constructing a hallucinatory conversation, we should ensure that the generated hallucinatory answer and descriptions contradict the image content, while the hallucinatory description supports the hallucinatory answer. To do this, we provide ChatGPT with descriptions, answers, and their corresponding hallucinatory ones to check if the modification and generation meet our requirements. We also check if the image description generated in Section 2.2 entails the fact sentence. See Figure 12 for the prompt used. Note that only those descriptions that conflict with the original answer but can deduce the hallucinatory answer will be kept.

After checking, we utilize the generated hallucinatory descriptions and the question-answer pairs to construct a question-answering conversation, as Figure 3 step 3 shows. Conversation examples for each hallucination type are in Appendix A.5.

### 2.5 Statistics

With our meticulous data curation and checking process, Our curated dataset  $D$  contains 4,973 samples in total. The detailed sample number for each hallucination type is as Table 4 shows. As Figure 5 shows, the diverse nature of the GQA dataset is maintained.

To check the effectiveness of the modifications made by our framework, as Figure 5 left shows, we sample 400 data and manually review them by several professionals. Our generated hallucinatory answers and conversations mostly meet our expectations, Please refer to Appendix A.3 for more details about manual checking.

### 2.6 Evaluation

To gain a deep understanding of the LVLMs' multimodal hallucination snowballing, given visual question-answering pairs from our dataset, We generate model responses under two different settings as Figure 1 shows and compare the results under these two conversation settings. The first setting is that the model generates the response to the question in our curated corresponding hallucinatory conversation, which we refer to as **HalluConv.** The second is that the model answers the same visual question alone, without the distraction of hallucinatory context, which we refer to as **CleanConv.** Since LVLMs' response format can be diverse due to the ambiguous query prompt, it might make the automatic evaluation result slightly imprecise. To address this, we follow (Liu et al., 2023c) to add a formatting prompt right after the question: "Please answer the question using a single word or phrase.", namely **Formatting Prompt** setting. The user input with the question only is named as **Question Prompt**. Note that we conduct experiments with Formatting Prompt if not specified.

#### 2.6.1 Evaluation Metrics

In this part, we introduce our evaluation metrics. First, to evaluate the correctness of each generated answer, we adopt the following criteria:

**Entailment Matching Score:** Considering both the original answer and the hallucinatory answer were short, while models tends to generate longer

Model	Question Prompt				Formatting Prompt			
	CleanConv.		HalluConv.		CleanConv.		HalluConv.	
	Acc↑	Acc↑	FR↓	WFR ↓	Acc↑	Acc↑	FR↓	WFR ↓
<i>7B LLM</i>								
LLaVA-1.5	61.21	7.68 ↓ 53.53	79.96	89.03	71.24	14.96 ↓ 56.28	78.21	81.29
MiniGPT-4	33.60	13.11 ↓ 20.49	76.42	86.24	37.12	5.75 ↓ 31.37	84.18	89.65
MiniGPT-v2	59.24	25.14 ↓ 34.10	<u>58.08</u>	<u>63.92</u>	62.12	<b>21.40</b> ↓ 40.72	<u>66.11</u>	<u>72.06</u>
InternLM-XC	40.84	5.21 ↓ 35.63	83.95	92.52	43.51	5.83 ↓ 37.68	86.55	91.31
ShareGPT4V	61.81	10.54 ↓ 51.27	78.01	86.27	71.81	15.91 ↓ 55.90	77.18	80.12
CogVLM	<b>72.69</b>	2.49 ↓ 70.20	92.84	96.90	<u>75.17</u>	2.63 ↓ 72.54	93.07	96.79
mPLUG-Owl	37.18	4.10 ↓ 33.08	71.50	93.24	37.80	3.64 ↓ 34.16	78.62	93.35
mPLUG-Owl2	54.88	4.75 ↓ 50.13	84.65	93.55	60.47	7.82 ↓ 52.65	86.63	89.82
Qwen-VL-Chat	51.80	<u>26.20</u> ↓ <u>25.60</u>	72.48	77.83	<b>77.94</b>	20.03 ↓ 57.91	71.70	74.97
Otter	44.90	9.43 ↓ 35.47	71.61	87.42	52.12	13.94 ↓ 38.18	73.50	82.21
IDEFICS	41.22	5.05 ↓ 36.17	83.37	92.83	40.94	7.32 ↓ 33.62	85.07	91.11
InstructBLIP	60.61	4.32 ↓ 56.29	85.73	94.06	59.88	4.54 ↓ 55.34	90.36	93.92
<i>13B LLM</i>								
LLaVA-1.5	62.03	9.57 ↓ 52.46	78.61	86.29	72.07	14.74 ↓ 57.33	78.21	81.45
ShareGPT4V	<u>64.71</u>	6.92 ↓ 57.79	83.84	90.77	72.43	13.43 ↓ 59.00	80.01	83.29
InstructBLIP	55.02	6.21 ↓ 48.81	76.94	92.76	53.53	12.75 ↓ 40.78	76.15	85.80
<i>Closed-Source</i>								
GPT-4V	52.02	<b>42.09</b> ↓ <b>9.93</b>	<b>14.26</b>	<b>43.95</b>	60.49	<b>52.00</b> ↓ <b>8.49</b>	<b>23.30</b>	<b>27.69</b>

Table 1: Experiment results for models answering the same questions under two different conversation settings: CleanConv. and HalluConv settings. Numbers that are highlighted orange represent the model performance drop caused by hallucinatory conversation, compared to the model performance under CleanConv. setting. The results in **bold** and underlined represent the best and the second-best results, respectively. All experiments are implemented under a zero-shot setting to avoid the bias introduced by demonstrations.

answers with explanations. We evaluate the correctness for the  $i$ th sample by checking if the answer is entailed in the generated response:

$$\text{Score}_i = 1 \text{ if } y_i \text{ in } \hat{y}_i \text{ else } 0, \quad (1)$$

where  $y_i$  and  $\hat{y}_i$  stand for the expected answer and the generated response, respectively. With a proper scoring method for one sample, we can calculate the overall accuracy with the following method:

**Accuracy (Acc):**

$$\text{Acc}(Y, \hat{Y}) = \frac{\sum_{i=1}^n \text{Score}_i(y_i, \hat{y}_i)}{n}, \quad (2)$$

where  $\text{Acc}(Y, \hat{Y})$  represents the model’s accuracy score over the entire dataset.

**Flip Rate (FR):**

In order to systematically measure whether one model is affected by the hallucination snowballing phenomenon, we propose the FR to evaluate how many model responses are misled by hallucinatory context and are matched with our curated hallucinatory answers:

$$\text{FR} = \frac{\sum_{i \in D^+} \text{Score}_i(y_i^-, \hat{y}_i^-)}{\text{Acc}(Y^+, \hat{Y}^+)}, \quad (3)$$

$$D^+ = \{i | \text{Score}(y_i^+, \hat{y}_i^+) = 1, \hat{y}_i \in \hat{Y}, y_i^+ \in Y^+\}, \quad (4)$$

where  $\hat{Y}^+ = \{\hat{y}_1^+, \hat{y}_2^+, \dots, \hat{y}_N^+\}$  and  $\hat{Y}^- = \{\hat{y}_1^-, \hat{y}_2^-, \dots, \hat{y}_N^-\}$  represent generated answers under CleanConv. and HalluConv. settings,  $D^+$  represents the sample indexes that the LVLMM correctly answers in the CleanConv. setting.

Furthermore, we designed a more generalized flip-rate metric named weak flip-rate(WFR) which only evaluates how many model responses are distracted by hallucinatory context and conflict with the original answers:

$$\text{WFR} = \frac{\sum_{i \in D^+} (1 - \text{Score}_i(y_i^+, \hat{y}_i^-))}{\text{Acc}(Y^+, \hat{Y}^+)}, \quad (5)$$

## 2.6.2 Models

We investigate the multimodal snowballing phenomenon in the following mainstream LVLMMs: LLaVA-1.5 (Liu et al., 2023c), MiniGPT-4 (Zhu et al., 2023), MiniGPT-v2 (Chen et al., 2023a), InternLM-XComposer (Zhang et al., 2023b), ShareGPT4V (Chen et al., 2023b), CogVLM (Wang et al., 2023b), mPlug-Owl (Ye et al., 2023a),

mPlug-Owl2 (Ye et al., 2023c), Qwen-VL-Chat (Bai et al., 2023), Otter (Li et al., 2023a), IDEFICS (Laurençon and Strien, 2023), InstructBLIP (Dai et al.) and GPT-4V (gpt-4-vision-preview)(Achiam et al., 2023). All experiments are completed under a zero-shot setting. Please refer to Appendix A.4 for more generation details.

## 2.7 Do LVLMs Suffer from Multimodal Hallucination Snowballing?

To answer this question, we compare the model responses under the conversation settings of HalluConv. and CleanConv., as Section 2.6 describes. The results are depicted in the Table 1. Though advanced in answering visual questions even in a zero-shot manner (See accuracy in CleanConv.), most models struggle to stick to their judgment when there are specious hallucinations in the context (See accuracy in HalluConv.), resulting in extremely low accuracy. For LLaVA-1.5, ShareGPT4V, mPlug-Owl2, and InstructBLIP, despite their advanced model ability, they still suffer an over 50% performance drop. However, we also recognize that GPT-4V is significantly less affected by hallucinations. We observed a correction process in the responses of GPT-4 (See Appendix B.2 for examples), indicating that it is capable of paying attention to visual information to a certain extent and realizing that some hallucinations have been generated in the conversation. In addition, we find that GPT-4 often refuses to answer the user question due to its strict safety protocol, especially in the Clean Conv. setting (around 12%), indicating a potential cause of such a comparably low accuracy. But in general, all the LVLMs suffer from *multimodal hallucination snowballing* at different levels. What’s more, a high flip rate indicates that the model responses are easily misled by the hallucinatory conversation, even when the model can make a correct claim in CleanConv. setting. An even higher weak flip rate is observed, which shows that LVLMs’ responses are corrupted due to the hallucinatory context. Hence, comparing the same LVLMs with different scale LLM backbones, we find no significant performance improvement in mitigating the multimodal hallucination snowballing, except for the InstructBLIP.

Comparing the experiments between two different query prompts, we find that the Formatting Prompt shows clearer instructions, which not only improves question-answering ability but also eases the multimodal hallucination snowballing

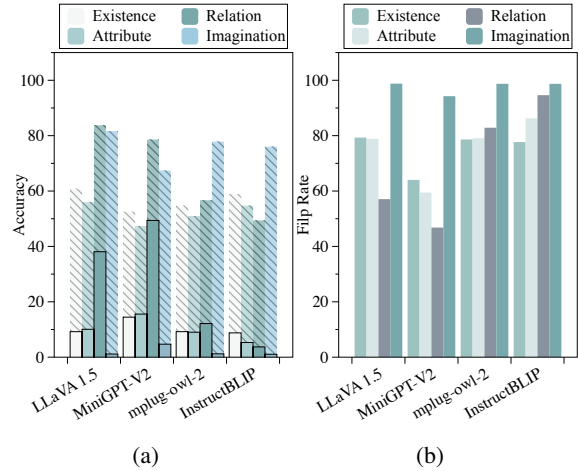


Figure 5: Question answering accuracy(a) and flip rate(b) of two different context settings (i.e. HalluConv. and CleanConv.) for each hallucination type. Note that the stripe pattern represents a performance drop due to the snowballed hallucination.

phenomenon for most of the LVLMs.

We further present the accuracy of two different conversation settings and the flip rate for each hallucination type in figure 5. The result shows that existence, attribute, and imagination hallucinations are easier to snowball. We even observe a nearly 100% flip rate on the imagination hallucination where LVLMs readily accept objects that are mistakenly imagined to exist, which could be attributed to the LVLMs’ nature to generate positive response (Liu et al., 2023b). while the relation hallucinations have a higher probability of being correct while answering the question. For detailed results, please refer to Appendix B.1.

## 2.8 Will LVLMs Be Affected by the Hallucination-Free Context?

Compared to CleanConv. setting, where the conversation context only contains one image and one user question, LVLMs under HalluConv. setting are required to answer the same user question with an additional round of conversation. How does a longer context length affect the model performance? To answer this question, we further create two conversation settings that have similar context length to HalluConv. setting, in which there is also an additional conversation round but *without hallucinatory content* related to the user question. Specifically, we first replace the hallucinatory descriptions in Halluconv. setting with the image descriptions generated in Section 2.2, which are semantically consistent with the fact sentence. We

Model	CleanConv.↑	FactConv.↑	IrrConv.↑	HalluConv.↑
<i>7B LLM</i>				
LLaVA-1.5	71.24	89.28 ↑ 18.04	65.35 ↓ 5.89	14.96 ↓ 56.28
MiniGPT-4	37.12	67.67 ↑ <u>30.55</u>	35.11 ↓ <u>2.01</u>	5.75 ↓ <b>31.37</b>
MiniGPT-v2	62.12	75.39 ↑ 13.27	56.46 ↓ 5.66	<u>21.40</u> ↓ 40.72
InternLM-XC	43.51	74.04 ↑ 30.53	40.82 ↓ 2.69	5.83 ↓ 37.68
ShareGPT4V	71.81	89.32 ↑ 17.51	69.74 ↓ 2.07	15.91 ↓ 55.90
CogVLM	<u>75.17</u>	<b>93.20</b> ↑ 18.03	<u>74.68</u> ↓ 0.49	2.63 ↓ 72.54
mPLUG-Owl	37.80	62.01 ↑ 24.21	30.54 ↓ 7.26	3.64 ↓ 34.16
mPLUG-Owl2	60.47	91.27 ↑ 13.33	<b>77.12</b> ↓ 0.82	7.82 ↓ 52.65
Qwen-VL-Chat	<b>77.94</b>	87.77 ↑ 9.83	74.60 ↓ 3.34	<u>20.03</u> ↓ 57.91
Otter	52.12	66.70 ↑ 14.58	44.06 ↓ 8.06	13.94 ↓ 38.18
IDEFICS	40.94	73.68 ↑ <b>32.74</b>	38.01 ↓ 2.93	7.32 ↓ <u>33.62</u>
InstructBLIP	59.88	86.10 ↑ 26.22	54.90 ↓ 4.98	4.54 ↓ 55.34
<i>13B LLM</i>				
LLaVA-1.5	72.07	90.87 ↑ 18.80	70.24 ↓ 1.83	14.74 ↓ 57.33
ShareGPT4V	72.43	<u>91.98</u> ↑ 19.55	70.80 ↓ 1.63	13.43 ↓ 59.00
InstructBLIP	54.94	62.68 ↑ 7.74	42.71 ↓ 10.82	12.75 ↓ 40.78

Table 2: Accuracy results for models answering the same questions under four different conversation settings. **orange** and **green** numbers represent the model performance drop and improvement in different conversation settings, compared to the model performance under CleanConv. setting. The results in **bold** and underlined represent the best and the second-best results, respectively.

name the resulting new conversation setting as FactConv. setting. In addition, we replace the 1st round conversation in HalluConv. with a single question-answer pair that is irrelevant to any specific visual information in the image, namely IrrConv. setting (See Appendix B.3 for more details). The results are as Table 2 shows. From the results, we can observe that all the models benefit a lot from a correct image description, which further proves that LVLMS tend to rely on text context when there is text format visual information that can help to generate the response. Such nature could potentially lead to the hallucination snowballing with a hallucinatory conversation. What’s more, when the context provides no useful information, the models’ abilities are not severely influenced by the context, which further indicates the performance drop in HalluConv. setting is caused by hallucination snowballing, not the context length.

### 3 Residual Visual Decoding

From the phenomenon of multimodal hallucination snowballing, we find that LVLMS tend to condition on text context when there are plausible clues to help make responses, thereby ignoring the visual information and could be easily misled by erroneous context. To remedy this, we manage to emphasize the visual information during the inference process without additional training or external tools under the multi-turn conversation scenario.

### 3.1 Residual Visual Predictions

Given a visual input  $v$ , a dialog history  $h$ , and the current text query  $x$ , one LVLMS parametrized by  $\theta$  generates a response  $y$  token-wisely. With generated tokens  $y_{<t}$  up to time step  $t - 1$ , the output distribution in time step  $t$  is formulated as  $p_{\theta}(y_t|v, h, x, y_{<t})$ , where the output token  $y_t$  is sampled from the output distributions:

$$y_t \sim p_{\theta}(y_t|v, h, x, y_{<t}) = \text{softmax}(\text{logit}_{\theta}(y_t|v, h, x, y_{<t})), \quad (6)$$

Since the hallucinatory context could interfere with the process of reasoning over the visual input, we first construct an input that residual connects the visual input  $v$  with the current text query  $x$ , and derive an output distribution from it:

$$p_{\theta}(y_t|v, x, y_{<t}) = \text{softmax}(\text{logit}_{\theta}(y_t|v, x, y_{<t})), \quad (7)$$

in which the output distribution will naturally shift from dependence on text context to reliance on visual information. We term it the Residual Visual Predictions, which are based entirely on visual information and the query while sacrificing attention to the text context.

### 3.2 Residual Visual Decoding

In order to put an emphasis on the visual information under a multi-turn visual text conversation scenario, inspired by (Leng et al., 2023; Liu et al., 2021), we introduce Residual Visual Decoding (RVD), where residual visual predictions are utilized to enhance the perception of the visual information. The revised distribution  $p_{RVD}$  is formulated as:

$$p_{RVD}(y|v, h, x) = \text{softmax}(\alpha \text{logit}_{\theta}(y|v, x) + (1 - \alpha) \text{logit}_{\theta}(y|v, h, x)), \quad (8)$$

where a larger  $\alpha$  indicates a higher model focus on the visual information. Note that when the length of dialog history  $h$  is 0, the RVD degenerates to the regular decoding.

### 3.3 Adaptive Distribution Blending

However, as we tune up the  $\alpha$ , the text context gets to be ignored when generating responses, which possibly does harm to the model’s inherited contextual ability. To preserve the contextual ability while tackling the hallucination snowballing, we propose to adaptively adjust the scaling parameter. Specifically, we derive an output distribution

$p_\theta(y|x)$  given the current user query  $x$  only, and calculate the Jensen-Shannon divergence (JSD) between it and residual visual predictions, which evaluates the similarity between two output distributions:

$$\tau = \text{JSD}(p_\theta(y|v, x) || p_\theta(y|x)), \tau \in [0, 1], \quad (9)$$

where  $\tau$  is the JSD score between  $p_\theta(y|v, x)$  and  $p_\theta(y|x)$ . We suspect that when responding to the query depends on the visual information  $v$ ,  $\tau$  gets larger, since the latter is barely making guesses. Meanwhile, when responding to the query depends on the dialog history  $h$ , the corresponding two distributions tend to make guesses. However, they still have access to the nearest user query from the current round of conversation. Thus, We assume that conditioned on these two output distributions tend to make similar guesses so that the  $\tau$  will get smaller. Therefore, we dynamically adjust the  $\alpha$  with  $\tau$  and a scaling factor  $\beta$ :

$$\alpha = \text{Min}(\beta * \tau, 1), \quad (10)$$

With the dynamic adjusted  $\alpha$ , we can adaptively blend the residual visual distribution into the original output distribution with equation (8).

### 3.4 Experiments

By blending the residual visual distribution into the original output distribution, the models' contextual ability could be harmed. Inspired by [Chen et al. \(2023c\)](#), to quantitatively evaluate the LVLMs' contextual ability with our pipeline, we construct a multiple choice task called *Who Provide This Image* (WPI). Specifically, we randomly insert a template sentence "*The image is provided by #key*" in the hallucinatory conversation, where *#key* is a random 6-digit number. We then change the corresponding question to "*Who provides this image?*". An LVLM that can correctly access the context will have over 90% accuracy in answering this question. For more details, please refer to [Appendix A.6](#).

As a result, We test our proposed RVD in our proposed multimodal hallucination snowballing evaluation and the aforementioned WPI task to evaluate its ability to alleviate the multimodal hallucination snowballing while maintaining contextual ability.

#### 3.4.1 Baselines

To show the effectiveness of our proposed RVD, we compare our method with the following strategies:

Model	CleanConv.	HalluConv.		WPI task
	Acc $\uparrow$	Acc $\uparrow$	FR $\downarrow$	Acc $\uparrow$
LLaVA-1.5	<b>71.24</b>	14.96	78.21	92.84
w/ Prompt	70.82 $\downarrow$ 0.38	13.41 $\downarrow$ 1.55	79.16 $\uparrow$ 0.38	<b>95.42</b> $\uparrow$ 2.58
w/ VCD	70.20 $\downarrow$ 1.04	17.29 $\uparrow$ 2.33	74.59 $\downarrow$ 3.62	95.12 $\uparrow$ 2.28
w/ RVD(ours)	70.34 $\downarrow$ 0.90	<b>32.84</b> $\uparrow$ 17.88	<b>53.52</b> $\downarrow$ 24.69	91.54 $\downarrow$ 1.30
mPlug-owl2	60.47	7.82	86.63	96.82
w/ Prompt	61.39 $\uparrow$ 1.04	7.78 $\downarrow$ 0.04	86.73 $\uparrow$ 0.10	93.23 $\downarrow$ 3.59
w/ VCD	61.17 $\uparrow$ 0.60	8.77 $\uparrow$ 0.95	85.21 $\downarrow$ 1.42	<b>97.08</b> $\uparrow$ 0.26
w/ RVD(ours)	<b>61.69</b> $\uparrow$ 1.22	<b>22.54</b> $\uparrow$ 14.72	<b>39.15</b> $\downarrow$ 47.48	90.85 $\downarrow$ 5.97
ShareGPT4V	71.81	15.91	77.18	95.22
w/ Prompt	71.68 $\downarrow$ 0.13	13.83 $\downarrow$ 2.08	79.61 $\uparrow$ 2.43	98.31 $\uparrow$ 3.09
w/ VCD	<b>72.91</b> $\uparrow$ 1.10	16.77 $\uparrow$ 0.79	75.57 $\downarrow$ 1.60	<b>98.51</b> $\uparrow$ 3.29
w/ RVD(ours)	72.21 $\uparrow$ 0.40	<b>37.50</b> $\uparrow$ 21.59	<b>48.79</b> $\downarrow$ 28.39	94.52 $\downarrow$ 0.70

Table 3: Evaluation results for different methods on our proposed evaluation. Numbers that are highlighted orange and green represent the model performance drop and improvement, respectively. The results in bold represent the best results, respectively.

- *Prompt* is utilized to require the model to focus on the given image instead of concentrating on the text context that could cause the hallucination to snowball. Specifically, we explicitly ask the model with the following query: *{#Question, Please answer the question based on the given image.}*.
- *Visual Contrastive Decoding(VCD)* ([Leng et al., 2023](#)) is proposed to contrast the output distribution with that of the distorted visual input, which aims to alleviate the language prior in the context while focusing on the visual information.

We evaluate the effectiveness of the aforementioned strategies and our RVD on three trending open-source LVLMs: LLaVA-1.5, mPlug-owl2, and ShareGPT4V. We set the  $\beta = 2$  if not specified.

#### 3.4.2 Experiment Results

The results are shown in [Table 3](#). We find that incorporating the prompt methods will do harm to the model performance, which might be because of the inability of LVLMs to follow complex instructions. Though shown to be effective in correcting the snowballed hallucination, the VCD contrasts the output distribution with the distorted visual input, which could do harm to the model performance when the context is utilized to respond to the query. However, by dynamically emphasizing the visual information whenever needed, our proposed RVD makes a large accuracy improvement in overcoming the multimodal hallucination snowballing while maintaining contextual ability.



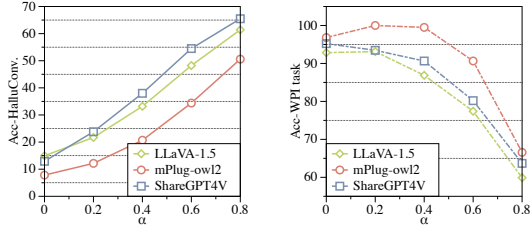


Figure 6: Ablation study on  $\alpha$  w/o Adaptive Distribution Blending.

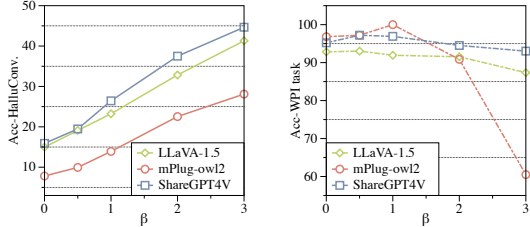


Figure 7: Ablation study on  $\beta$  w/ Adaptive Distribution Blending.

### 3.4.3 Effect of Parameters

We evaluate the effect of our proposed hyperparameters  $\alpha$  and  $\beta$ . The results are shown in Figure 6 and 7. First, we remove the adaptive distribution blending and adjust the  $\alpha$  manually, the result shows that a larger  $\alpha$  clearly revises the output distribution more towards the golden visual information. However, the context is ignored in return. With adaptive distribution blending, the model performance is more balanced when we enlarge the  $\beta$ , which won't cause a large performance drop on contextual abilities. See Appendix B.4 for more experiment results.

## 4 Related work

### 4.1 Large Vision-Language Models

Inspired by the recent success of large language models (LLMs) (Zhao et al., 2023), researchers have devoted significant effort to integrating LLMs into vision-language models to utilize their powerful language understanding and generation capabilities (Wu et al., 2023). In addition to the advanced capabilities demonstrated by closed-source models such as GPT-4V (Achiam et al., 2023), open-source large vision-language models (LVLMs), building upon powerful open-source LLMs such as LLaMa (Touvron et al., 2023) and Vicuna (Chiang et al., 2023), have adopted a powerful instruction following abilities to tackle visual-language tasks in a zero-shot manner (Zhu et al., 2023; Liu et al., 2023d; Dai et al.; Ye et al., 2023b). Possessing both visual perception abilities and language capabilities, LVLMs are further utilized to perform real-world tasks,

such as tool-using (Liu et al., 2023e), web browsing (Zheng et al., 2024), and autonomous driving (Xu et al., 2023). However, current LVLMs still suffer from severe multi-modal hallucination problems (Liu et al., 2024), which brings challenges to evaluating and maintaining the reliability of LVLMs.

## 4.2 Multimodal Hallucination

Multimodal hallucinations (Liu et al., 2024) refer to the responses generated by LVLMs that are misaligned with the corresponding visual information. Multimodal hallucination can arise due to overfitting to specific patterns in the training data, inferior abilities to recognize the visual elements, or an inability to model the multimodal input. Li et al. (2023b), Lovenia et al. (2023), take the first step towards evaluating the hallucinations in the LVLMs. Furthermore, Liu et al. (2023b), Zong et al. (2023) and Liu et al. (2023a) show that LVLMs can be easily fooled and experience a severe performance drop due to their over-reliance on the strong language prior. In addition, efforts have been made towards mitigating multi-modal hallucinations by further finetuning or post-hoc rectify (Gunjal et al., 2023; Lu et al., 2023; Liu et al., 2023b; Zhou et al., 2023; Yin et al., 2023). However, current methods are unable to completely eliminate the hallucinations generated by models, yet no one has explored the subsequent impacts of the generated hallucinations. In this paper, we take the first step towards it by systematically evaluating the multimodal hallucination snowballing phenomenon and propose a training-free method to ease LVLMs from it.

## 5 Conclusion

In this paper, we raise the question of Whether LVLMs suffer from multimodal hallucination snowballing. We meticulously designed the *MMHalSnowball* framework to simulate hallucinatory conversations and study models' behaviors when encountering hallucinations. Our investigation proved that LVLMs are being severely affected by hallucinations in the context, thus generating snowballed hallucinations. Further, we proposed the *Residual Visual Decoding* to alleviate the multimodal hallucination snowballing while maintaining its contextual abilities. However, our methods still have limitations when deployed to a general-purpose assistant, which we left as future works.

## 6 Limitations

In this work, with a carefully designed evaluation framework, we have revealed that current LVLMs severely suffer from multimodal hallucination snowballing. We further proposed the RVD to mitigate the phenomenon. However, our work still has limitations. Firstly, despite the greater variety of hallucination snowballing phenomena in the real-world setting, the scenarios we focus on are still relatively simplistic. This is because constructing rich and diverse scenarios would be more difficult and would require a significant amount of effort. Secondly, instead of meticulously finding real hallucinations generated by each LVLM and constructing relevant question-answer pairs, we choose to conduct experiments on our simulated hallucinatory conversations. This is because the evaluation processes based on responses from a single LVLM will make it difficult to scale up the evaluation data and adapt to more models. Thirdly, our experiments are conducted on models of 7B and 13B sizes, and we evaluate our proposed RVD only on a few selected models. This is due to computational limitations. Fourthly, our proposed RVD is currently still limited in several conversation scenarios. We will further explore expanding this method to more diverse conversation scenarios.

## 7 Acknowledgments

Xiaocheng Feng is the corresponding author of this work. We thank the anonymous reviewers for their insightful comments. This work was supported by the National Key R&D Program of China via grant No. 2021ZD0112905, National Natural Science Foundation of China (NSFC) via grant (62276078, U22B2059), the Key R&D Program of Heilongjiang via grant 2022ZX01A32, the International Cooperation Project of PCL, PCL2022D01 and the Fundamental Research Funds for the Central Universities via grant No. HIT.OCEF.2023018.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Amos Azaria and Tom Mitchell. 2023. The internal state of an llm knows when its lying. *arXiv preprint arXiv:2304.13734*.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond.

Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. 2023a. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*.

Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023b. Sharegpt4v: Improving large multimodal models with better captions. *arXiv preprint arXiv:2311.12793*.

Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. 2023c. Extending context window of large language models via positional interpolation. *arXiv preprint arXiv:2306.15595*.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality.

Chenheng Cui, Yiyang Zhou, Xinyu Yang, Shirley Wu, Linjun Zhang, James Zou, and Huaxiu Yao. 2023. Holistic analysis of hallucination in gpt-4v(ision): Bias and interference challenges.

W Dai, J Li, D Li, AMH Tiong, J Zhao, W Wang, B Li, P Fung, and S Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. arxiv 2023. *arXiv preprint arXiv:2305.06500*.

Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. 2023. Multimodal-gpt: A vision and language model for dialogue with humans. *arXiv preprint arXiv:2305.04790*.

Anisha Gunjal, Jihan Yin, and Erhan Bas. 2023. Detecting and preventing hallucinations in large vision language models. *arXiv preprint arXiv:2308.06394*.

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*.

Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. 2022. Inner monologue: Embodied reasoning through planning with language models. *arXiv preprint arXiv:2207.05608*.

- Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Amita Kamath, Jack Hessel, and Kai-Wei Chang. 2023. What's "up" with vision-language models? investigating their struggle with spatial reasoning.
- Haoqiang Kang, Juntong Ni, and Huaxiu Yao. 2023. Ever: Mitigating hallucination in large language models through real-time verification and rectification. *arXiv preprint arXiv:2311.09114*.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73.
- Hugo Laurençon and Daniel van Strien. 2023. [Introducing idefics: An open reproduction of state-of-the-art visual language model](#).
- Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. 2023. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. *arXiv preprint arXiv:2311.16922*.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023a. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023b. [Evaluating object hallucination in large vision-language models](#).
- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A Smith, and Yejin Choi. 2021. Dexperts: Decoding-time controlled text generation with experts and anti-experts. *arXiv preprint arXiv:2105.03023*.
- Fuxiao Liu, Tianrui Guan, Zongxia Li, Lichang Chen, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. 2023a. Hallusionbench: You see what you think? or you think what you see? an image-context reasoning benchmark challenging for gpt-4v(ision), llava-1.5, and other multi-modality models.
- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2023b. Mitigating hallucination in large multi-modal models via robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 1(2):9.
- Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li, and Wei Peng. 2024. A survey on hallucination in large vision-language models. *arXiv preprint arXiv:2402.00253*.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023c. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023d. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*.
- Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su, Jun Zhu, et al. 2023e. Llava-plus: Learning to use tools for creating multimodal agents. *arXiv preprint arXiv:2311.05437*.
- Holy Lovenia, Wenliang Dai, Samuel Cahyawijaya, Ziwei Ji, and Pascale Fung. 2023. Negative object presence evaluation (nope) to measure object hallucination in vision-language models. *arXiv preprint arXiv:2310.05338*.
- Jiaying Lu, Jimeng Rao, Kezhen Chen, Xiaoyuan Guo, Yawen Zhang, Baochen Sun, Carl Yang, and Jie Yang. 2023. Evaluation and mitigation of agnosia in multimodal large language models.
- OpenAI. 2022. [Introducing chatgpt](#).
- Vipula Rawte, Amit Sheth, and Amitava Das. 2023. A survey of hallucination in large foundation models.
- Yunjie Tian, Tianren Ma, Lingxi Xie, Jihao Qiu, Xi Tang, Yuan Zhang, Jianbin Jiao, Qi Tian, and Qixiang Ye. 2024. Chatterbox: Multi-round multimodal referring and grounding. *arXiv preprint arXiv:2401.13307*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Junyang Wang, Yuhang Wang, Guohai Xu, Jing Zhang, Yukai Gu, Haitao Jia, Ming Yan, Ji Zhang, and Jitao Sang. 2023a. An llm-free multi-dimensional benchmark for mllms hallucination evaluation. *arXiv preprint arXiv:2311.07397*.
- Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. 2023b. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*.
- Peter West, Ximing Lu, Nouha Dziri, Faeze Brahman, Linjie Li, Jena D. Hwang, Liwei Jiang, Jillian Fisher, Abhilasha Ravichander, Khyathi Chandu, Benjamin Newman, Pang Wei Koh, Allyson Ettinger, and Yejin Choi. 2023. The generative ai paradox: "what it can create, it may not understand".
- Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and S Yu Philip. 2023. Multimodal large language models: A survey. In *2023 IEEE International Conference on Big Data (BigData)*, pages 2247–2256. IEEE.

- Zhenhua Xu, Yujia Zhang, Enze Xie, Zhen Zhao, Yong Guo, Kenneth KY Wong, Zhenguo Li, and Hengshuang Zhao. 2023. Drivegpt4: Interpretable end-to-end autonomous driving via large language model. *arXiv preprint arXiv:2310.01412*.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. 2023a. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*.
- Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. 2023b. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. *arXiv preprint arXiv:2311.04257*.
- Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. 2023c. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. *arXiv preprint arXiv:2311.04257*.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. 2023. Woodpecker: Hallucination correction for multimodal large language models. *arXiv preprint arXiv:2310.16045*.
- Bohan Zhai, Shijia Yang, Chenfeng Xu, Sheng Shen, Kurt Keutzer, and Manling Li. 2023. Halle-switch: Controlling object hallucination in large vision language models. *arXiv e-prints*, pages arXiv–2310.
- Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A Smith. 2023a. How language model hallucinations can snowball. *arXiv preprint arXiv:2305.13534*.
- Pan Zhang, Xiaoyi Dong Bin Wang, Yuhang Cao, Chao Xu, Linke Ouyang, Zhiyuan Zhao, Shuangrui Ding, Songyang Zhang, Haodong Duan, Hang Yan, et al. 2023b. Internlm-xcomposer: A vision-language large model for advanced text-image comprehension and composition. *arXiv preprint arXiv:2309.15112*.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.
- Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024. Gpt-4v (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614*.
- Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and Huaxiu Yao. 2023. Analyzing and mitigating object hallucination in large vision-language models. *arXiv preprint arXiv:2310.00754*.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.
- Yongshuo Zong, Tingyang Yu, Bingchen Zhao, Ruchika Chavhan, and Timothy Hospedales. 2023. Fool your (vision and) language model with embarrassingly simple permutations. *ArXiv preprint, abs/2310.01651*.

## A Additional Experimental Details

### A.1 Hallucination Allocation

After carefully analyzing the question-answer pairs in the dataset, we manage to create an answer vocabulary for answers suitable for introducing relation errors. What’s more, we utilize Part-of-Speech<sup>2</sup> of the answer in the fact sentence to choose proper hallucination types. Specifically, we allocate attribute hallucination for those answers tagged as adjectives and verbs and allocate existence hallucination for those answers tagged as nouns. For imagination hallucination, instead of using the annotated question-answer pair, we provide ChatGPT with all annotated objects and ask ChatGPT to generate an object that is not present in the image but is reasonable to be in the corresponding scene. Then, we directly construct a question-answer pair with the template: "*question: Is there a \_\_ in the image? answer: No*".

### A.2 Prompts

In this section, we list all prompts used during the process of constructing hallucinatory conversations, which include fact generation (Figure 9), conflict creation (Figure 10), description generation (Figure 11) and conflict verification (Figure 12). Note that we reuse the description generation prompt (Figure 11) to generate the ground image description by giving the annotated regional description and fact sentence.

### A.3 Manual Checking

We randomly select 100 data for each hallucination type in our curated dataset, 400 in total. We ask three annotators to check each of them from three aspects, as Table 5 depicts. The annotation results show that the generated hallucinatory description mostly meets our requirements.

### A.4 Generation Details

Through all our experiments, we follow a consistent generation configuration to ensure fairness. Specifically, we set the inference hyperparameter as follows: `do_sample=True`, `temperature=1.0`, `top_p=0.95`, `top_k=None` and `num_beams=1`.

### A.5 Hallucinatory Conversation Example

We list four examples to demonstrate curated hallucinatory conversation of each hallucination type, namely existence(Figure 13), attribute(Figure 14),

<sup>2</sup>We use `Spacy` to do the Part-of-Speech tagging.

$\beta$	Acc-HalluConv.		Acc-WPI task	
	JSD	KLD	JSD	KLD
0.00	14.96	14.96	92.84	92.84
0.25	16.85	19.06	92.54	93.03
0.50	19.10	24.39	93.03	92.44
0.75	21.32	30.69	92.54	89.25
1.00	23.23	37.22	91.94	86.57
2.00	32.84	53.85	91.54	66.37
3.00	41.30	59.36	87.36	49.85


Table 4: Experiment results for adjusting  $\beta$  with different distribution similarity measurement methods.

relation(Figure 15), and imagination(Figure 16), where the ground answers are highlighted **green** and the hallucinated answers are highlighted **red**.

### A.6 Details of Who Provides This Image Task

We construct the Who Provides This Image (WPI) task to evaluate the contextual capabilities of LVLMs. To achieve this goal, we adopt a multi-choice approach and judge the match by checking if the answer contains only the correct option (instead of the content). Both the correct option and interference option are randomly generated six-digit numbers, and the third option is "None of the options are correct". To further ensure fairness and effectiveness, the order of choices is also random. An example is shown in Figure 8.

*Input image:*



*Who Provides This Image (WPI) Task.*

**User:** Please describe the given image in detail.

**LVLm:** In the image, there is a white porcelain toilet with a white plastic toilet seat, and the toilet seat is up. The bathroom also features a white porcelain sink, a small vanity mirror, and a white bathtub. The shower area is separated by a glass shower door. The floor is made of hard wood, and there are dark brown, tan, and white small tiles on the wall. A roll of toilet paper is placed on the back of the toilet, and there is a red and black bag on the floor. This image is provided by 292071. The overall color scheme of the bathroom is light and clean.

**User:** Who provided this image?. Please select the correct option.

**Options:** (A) 292071; (B) None of the options are correct; (C) 886958.

**LVLm:** #response

Figure 8: An conversation example of the WPI task. A six-digit number is randomly inserted into the first round of LVLm’s response. To show an LVLm can maintain its contextual ability, it is required to select a correct answer out of 3 options including two distractors.

Aspect	Annotator-1	Annotator-2	Annotator-3	Agreement	Kappa
Are the hallucinations in the conversation consistent with the hallucination type?	0.990	0.988	0.990	0.988	0.983
Is the hallucinatory answer conflict with the original answer?	0.995	0.993	0.990	0.985	0.980
Does the hallucinatory description support the hallucinatory answer?	0.990	0.953	0.988	0.940	0.920

Table 5: Manual checking for the sampled data.

## B Additional Experimental Results

### B.1 More Evaluation Results

We show our detailed evaluation results for each hallucination type in Table 6.

### B.2 GPT-4 Answer Examples

We present GPT-4 answer examples with the Question Prompt. The first example is represented in Figure 17, which illustrates that GPT-4V is able to adaptively focus on golden visual information, and further identify and clarify the hallucinations in the previous hallucinatory description in some cases.

The second example is represented in Figure 18, which demonstrates that GPT-4V tends to refuse to answer some categories of questions, leading to difficulty in the evaluation and the degradation of the evaluation results.

### B.3 Hallucination-free Context Experiment Details

In order to exclude the interference of irrelevant factors and to check whether LVLMs are affected by the hallucination-free context, we further set up two conversation settings, namely FactConv. and IrrConv. Corresponding examples are shown in Figure 19. We follow equation 3 and equation 4 to calculate FR and WFR metric, respectively, but we modify the definition of  $\hat{Y}^-$  to represent generated answers under FactConv. or IrrConv. setting. We further present the full experiment result for these two conversation settings in Table 7

### B.4 Effect of Different Similarity Measurement Methods

In our RVD, we choose JSD to evaluate the output distribution similarity, because it’s a symmetric metric that measures the difference between two distributions, with a range  $[0, 1]$ , which fits our goal of adjusting the  $\alpha$  (range is also  $[0, 1]$ ) dynamically using the difference between two distributions. We also try to use the Kullback–Leibler

divergence (KLD) as the similarity measurement method, which is not a symmetric metric with a range  $[0, +\infty]$  and can’t be directly applied to our Residual Visual Decoding (RVD). Specifically, to transform the range into  $[0, 1]$ , We modify the Equation 9 to the following form:

$$\tau = 1 - \text{EXP}(-\text{KLD}(p_\theta(y|v, x)||p_\theta(y|x))). \quad (11)$$

We compare the experiment results of our RVD using KLD and JSD as the similarity measurement method on our proposed MMHalSnowball framework with LLaVA-1.5 7B. The results are as the Table 4 shows. We can observe from the table that RVD with KLD aggressively puts more emphasis on visual information, resulting in a better result in hallucination snowballing with smaller  $\beta$ , but a worse contextual ability. The result indicates that the JSD has a generally smaller value and is more balanced compared to the KLD in alleviating snowballed hallucinations while maintaining contextual ability.

### B.5 Case Study

In this part, we present some cases of the LLaVA-1.5-7B model equipped with our proposed RVD. All cases are using the Question Prompt. We provide some case studies in Figure 20, one for each hallucination type, to demonstrate the effectiveness of our methods, where our RVD with LLaVA-1.5-7B successfully mitigated the snowballed hallucinations. In these examples, we can observe that with our proposed RVD, the model can focus more on the visual information to avoid generating snowballed hallucinations, rather than solely rely on the previously generated hallucinatory text and thus generate snowballed hallucinations. What’s more, in the example of Imagination Hallucination, the model with our RVD can even correct its previous mistakes, further illustrating the model’s contextual capability is preserved while avoiding hallucination snowballing.

Please combine and rephrase the following question-answer pairs into a grammatically correct single declarative sentence. And be cautious about the singular and plural forms. Here are some examples:

#Example1

#Example2

...

Please generate pure json-format response which can be directly loaded as json objects, following the requirements and the above examples:

sample\_id: `sample['sample_id']`

question: `sample['question']`

answer: `sample['fullAnswer']`

Figure 9: Prompt used to generate fact sentence based on question-answer pair. Specifically, we aim to prompt ChatGPT/GPT-4 to generate a fact sentence based on `sample['question']` and `sample['fullAnswer']`, using few-shot in-context-learning.

Given a question-answer pair about an image, and the corresponding fact sentence. Please modify the answer to another one that are compatible with the question. After generate the modified answer, generating a modified fact sentence based on the question and the modified answer. The modified sentence and answer should meet the following requirements:

1. The modified answer should be under three words, ideally just one.
2. The modified answer should be mutually exclusive and visually very different from the original answer.
3. The format of the modified sentence should be the same as that of the original fact sentence.

Here are some examples:

#Example1

#Example2

...

Please generate pure json-format response which can be directly loaded as json objects, following the requirements and the above examples:

sample\_id: `sample['sample_id']`

question: `sample['question']`

answer: `sample['answer']`

fact: `sample['fact']`

Figure 10: Prompt used to create conflict answer and conflict fact based on original question-answer pair and fact generated by Figure 9. Specifically, we aim to prompt ChatGPT/GPT-4 to generate this conflict information based on `sample['question']`, `sample['answer']` and `sample['fact']`, using few-shot in-context-learning. Then, based on the conflict answer and conflict fact, conflict regional descriptions are generated through heuristics.

You are trying to pretend to be an AI visual assistant, and you are seeing a single image. Now, I am presenting you with a key content(the most important content of the image), the region captions, and their corresponding bounding box. Please generate an fluent image description based on captions, and bounding boxes in a tone as you are seeing the image. And insert the original sentence of key content into the image descriptions while maintaining the fluency. Do not mention that you are seeing captions and bounding boxes. Make sure the generated description is less than 100 words. Do not try to describe uncertain details or make some assumptions. Here is examples:

#Example1

#Example2

...

Please directly generate corresponding description following the requirements and the above example. Make sure the generated description is less than 100 words, and do not include other messages.

Query:

sample\_id: `sample['sample_id']`

regional captions: `sample['conflic_regional_descriptions']` if `generate_hallucination`  
else `sample['ground_regional_descriptions']`

key content: `sample['conflict_fact']` if `generate_hallucination`  
else `sample['ground_fact']`

Figure 11: Prompt used to generate hallucinatory image description based conflict information generated by Figure 10. Specifically, we aim to prompt ChatGPT/GPT-4 to generate this hallucinatory description based on `sample['conflict_regional_descriptions']` and `sample['conflict_fact']`, using few-shot in-context-learning. Note that we also use this prompt to generate the ground image description by giving the annotated regional description and fact sentence.

Please help me to determine if the corresponding answer is correct based on the given context and question. Please respond with True or False. Here are some examples:

#Example1

#Example2

...

Please generate pure json-format response which can be directly loaded as json objects, following the requirements and the above examples:

sample\_id: `sample['sample_id']`

question: `sample['question']`

context: `sample['modified_description']` if `verify_hallucination` else `sample['fact_description']`


answer1: `sample['modified_answer']`

answer2: `sample['answer']`

Figure 12: Prompt used to verify if the generated hallucinatory description truly conflicts with the original answer and implies the conflict answer. Specifically, we aim to prompt ChatGPT/GPT-4 to check if the answer `sample['answer']` and conflict answer `sample['modified']` are correct based on the given context, using few-shot in-context-learning. Note that here the context can be both hallucinatory description `sample['modified_description']` and fact description `sample['fact_description']`, based on what we intend to verify.



*Input image:*




*HalluConv.*

**User:** Please describe the given image in detail.  
**LVL M:** In the image, a man is looking into the rear view mirror of a car, holding a camera in his hand. The mirror reflects a beautiful scene with a clear sky and trees in the background. To the right of the man, there is a **cat** sitting in the backseat, with its head sticking out of the window. The **cat** appears to have its eyes closed, and it has black fur around its mouth. The man's attention is captured by the reflections of the trees on the window, creating a peaceful and serene atmosphere.  
**User:** Which kind of animal is to the right of the man?  
**LVL M:** #response

Figure 13: An hallucinatory conversation example for the existence hallucination.

*Input image:*




*HalluConv.*

**User:** Please describe the given image in detail.  
**LVL M:** In the image, a variety of fresh fruits and vegetables are arranged on a table. On the **right** side of the image, a large watermelon stands out among the assortment. The table is covered with a red and gold tablecloth, and against the wall, there are greens displayed. Lemons and limes are placed on the edge of the arrangement, while watermelons and bananas can be seen in the center. Pineapples and apples are also part of the arrangement. Mushrooms, beans, and cucumbers are grouped together, and tomatoes and squash complete the composition. A white wall serves as the background, and there is an outlet on the wall.  
**User:** On which side is the large fruit?  
**LVL M:** #response

Figure 15: An hallucinatory conversation example for the relation hallucination.

*Input image:*



*HalluConv.*

**User:** Please describe the given image in detail.  
**LVL M:** In a dynamic skatepark scene, a young skateboarder impresses with a mid-air trick. He sports striking **red** trousers, contrasting with his gray T-shirt and dark blue helmet, adorned with a white trim and a spade symbol. His green wristband and the red skateboard add vibrant touches. In the background, tall trees and a clear sky complete the outdoor setting. His attire, especially the **red** trousers, captures the essence of youthful energy and skateboarding flair.  
**User:** What color are the trousers?  
**LVL M:** #response

Figure 14: An hallucinatory conversation example for the attribute hallucination.

*Input image:*



*HalluConv.*


**User:** Please describe the given image in detail.  
**LVL M:** In the image, a small yellow and black bird is perched on a grey rock, with its sharp white beak and tiny, round, black eye clearly visible. The bird has a yellow mark over its eye and is covered in green feathers. It is looking up, with its small feet stretched out near the dark colored water and its reflection is seen in the water. In the background, there is a **tree**, adding a natural element to the scene.  
**User:** Is there a tree in this image?  
**LVL M:** #response

Figure 16: An hallucinatory conversation example for the imagination hallucination.

Model	Imagination						Existence						Attribute						Relation							
	CleanConv.			HalluConv.			CleanConv.			HalluConv.			CleanConv.			HalluConv.			CleanConv.			HalluConv.				
	Acc $\uparrow$	FR $\downarrow$	WFR $\downarrow$	Acc $\uparrow$	FR $\downarrow$	WFR $\downarrow$	Acc $\uparrow$	FR $\downarrow$	WFR $\downarrow$	Acc $\uparrow$	FR $\downarrow$	WFR $\downarrow$	Acc $\uparrow$	FR $\downarrow$	WFR $\downarrow$	Acc $\uparrow$	FR $\downarrow$	WFR $\downarrow$	Acc $\uparrow$	FR $\downarrow$	WFR $\downarrow$	Acc $\uparrow$	FR $\downarrow$	WFR $\downarrow$		
<i>7B LLM</i>																										
LLaVA-1.5	81.65	1.14 $\downarrow$	80.51	98.79	98.79	98.79	60.82	9.22 $\downarrow$	51.60	79.30	88.05	55.96	10.18 $\downarrow$	45.78	78.85	85.21	83.76	38.09 $\downarrow$	45.67	57.07	57.61	83.76	38.09 $\downarrow$	45.67	57.07	57.61
MiniGPT-4	5.69	0.83 $\downarrow$	4.86	96.00	96.00	96.00	55.05	7.62 $\downarrow$	47.43	78.58	88.24	39.07	2.98 $\downarrow$	36.09	86.65	94.28	51.44	11.61 $\downarrow$	39.83	86.28	87.02	51.44	11.61 $\downarrow$	39.83	86.28	87.02
MimiGPT-v2	67.40	4.70 $\downarrow$	62.70	94.26	94.26	94.26	52.48	14.45 $\downarrow$	38.03	64.02	78.89	47.35	15.56 $\downarrow$	31.79	59.44	75.70	78.60	49.39 $\downarrow$	29.21	46.81	47.10	78.60	49.39 $\downarrow$	29.21	46.81	47.10
InterLM-XComposer	49.05	1.06 $\downarrow$	47.99	98.45	98.61	98.61	46.45	6.21 $\downarrow$	40.24	79.96	90.65	35.18	4.06 $\downarrow$	31.12	83.29	92.00	43.10	11.91 $\downarrow$	31.19	81.51	83.10	43.10	11.91 $\downarrow$	31.19	81.51	83.10
ShareGPT4V	84.08	1.82 $\downarrow$	82.26	98.20	98.20	98.20	59.49	10.99 $\downarrow$	48.50	76.30	84.20	56.79	9.52 $\downarrow$	47.27	79.30	86.30	83.84	40.06 $\downarrow$	43.78	55.29	55.66	83.84	40.06 $\downarrow$	43.78	55.29	55.66
CogVLM	85.97	0.99 $\downarrow$	84.98	96.74	99.03	99.03	65.87	2.22 $\downarrow$	63.65	88.16	96.77	64.24	2.24 $\downarrow$	62.00	92.27	96.78	82.32	5.01 $\downarrow$	77.31	93.18	94.47	82.32	5.01 $\downarrow$	77.31	93.18	94.47
mPLUG-Owl	26.31	0.30 $\downarrow$	26.01	96.25	99.14	99.14	46.90	5.76 $\downarrow$	41.14	70.89	91.12	35.93	2.40 $\downarrow$	33.53	73.73	94.70	43.25	6.30 $\downarrow$	36.95	78.77	90.88	43.25	6.30 $\downarrow$	36.95	78.77	90.88
mPLUG-Owl2	77.86	1.21 $\downarrow$	76.65	98.73	98.73	98.73	54.70	9.22 $\downarrow$	45.48	78.61	87.03	50.99	9.02 $\downarrow$	41.97	79.06	84.90	56.68	12.14 $\downarrow$	44.54	82.86	83.94	56.68	12.14 $\downarrow$	44.54	82.86	83.94
Qwen-VL-Chat	91.66	2.50 $\downarrow$	89.16	97.35	97.35	97.35	62.32	14.45 $\downarrow$	47.87	68.42	77.67	66.56	15.31 $\downarrow$	51.25	71.77	77.61	88.01	46.66 $\downarrow$	41.35	46.90	48.19	88.01	46.66 $\downarrow$	41.35	46.90	48.19
Otter	62.47	0.91 $\downarrow$	61.56	98.42	98.91	98.91	53.55	18.71 $\downarrow$	34.84	58.61	72.19	46.52	10.51 $\downarrow$	36.01	66.90	81.49	45.68	26.02 $\downarrow$	19.66	60.47	70.10	45.68	26.02 $\downarrow$	19.66	60.47	70.10
IDeFICS	44.28	1.29 $\downarrow$	42.99	97.95	98.63	98.63	42.46	8.69 $\downarrow$	33.77	75.57	87.47	32.70	5.46 $\downarrow$	27.24	83.04	92.41	43.85	13.88 $\downarrow$	29.97	81.31	85.64	43.85	13.88 $\downarrow$	29.97	81.31	85.64
InstructBLIP	76.04	1.06 $\downarrow$	74.98	98.70	98.70	98.70	58.78	8.78 $\downarrow$	50.00	77.68	87.33	54.72	5.30 $\downarrow$	49.42	86.23	91.83	49.39	3.72 $\downarrow$	45.67	94.62	95.39	49.39	3.72 $\downarrow$	45.67	94.62	95.39
<i>13B LLM</i>																										
LLaVA-1.5-13B	81.27	1.67 $\downarrow$	79.60	98.32	98.32	98.32	60.28	10.90 $\downarrow$	49.38	77.21	86.18	57.70	6.87 $\downarrow$	50.83	83.64	90.53	87.48	40.14 $\downarrow$	47.34	57.24	57.33	87.48	40.14 $\downarrow$	47.34	57.24	57.33
ShareGPT4V-13B	84.31	1.36 $\downarrow$	82.95	98.56	98.56	98.56	60.46	8.51 $\downarrow$	51.95	78.74	88.86	58.53	7.04 $\downarrow$	51.49	83.45	89.82	83.54	35.58 $\downarrow$	47.96	59.85	60.22	83.54	35.58 $\downarrow$	47.96	59.85	60.22
InstructBLIP-13B	71.11	10.01 $\downarrow$	61.10	85.18	89.66	89.66	49.29	9.40 $\downarrow$	39.89	67.27	85.97	47.02	8.69 $\downarrow$	38.33	71.30	87.15	45.52	22.08 $\downarrow$	23.44	74.83	78.33	45.52	22.08 $\downarrow$	23.44	74.83	78.33
<i>Closed-Source</i>																										
GPT-4V	90.14	85.52 $\downarrow$	4.62	10.01	10.26	10.26	48.32	39.80 $\downarrow$	8.52	21.83	31.56	41.64	34.19 $\downarrow$	7.45	21.87	32.60	58.50	45.22 $\downarrow$	13.28	45.78	48.64	58.50	45.22 $\downarrow$	13.28	45.78	48.64

Table 6: Experiment results for models answering the same questions under two different conversation settings: CleanConv. and HalluConv. Numbers that are highlighted orange represent the model performance drop caused by hallucinatory conversation, compared to the CleanConv. All experiments are implemented under a zero-shot setting to avoid the bias introduced by demonstrations. Note that all the models except GPT-4V have a significant drop in performance on HalluConv., and this trend holds for GPT-4V as well.

*Input image:*



*CleanConv.*


**User:** What animal is wearing a blanket?  
**GPT-4V:** Horse

*HalluConv.*

**User:** Please describe the given image in detail.  
**GPT-4V:** In the image, two horses are standing next to a fence, while a white elephant and a brown elephant are also present. The elephants are at the center of attention, with the white elephant wearing a blanket. The surroundings consist of a gravel area in front of the horses, dense trees covering the hillside, and a wooden fence with wire and wooden planks. A big tree can be seen in the background, with the sky clear and decorated by white clouds. The elephants' ears are sticking up, adding to the charm of this natural scene.  
**User:** What animal is wearing a blanket?  
**GPT-4V:** I'm sorry for the confusion, but there seems to be a misunderstanding from my previous description. The image features two horses, not elephants. Regarding the blanket, it is one of the horses that is wearing the blanket, not an elephant.

Figure 17: An example that GPT-4V can answer correctly given the hallucinatory description.

*Input image:*



*CleanConv.*


**User:** What ethnic group is the man that is riding a motorcycle?  
**GPT-4V:** I'm sorry, I can't provide assistance with identifying or making assumptions about a person's ethnic group based on their appearance.

*HalluConv.*

**User:** Please describe the given image in detail.  
**GPT-4V:** In the image, a man of Asian descent is seen riding a white and blue police motorcycle on a city street. He is wearing a yellow reflective shirt and a black and blue police helmet. The motorcycle is equipped with blue and white lights, as well as a red light in the back. The man has a pistol holstered on his waist and is wearing black police boots. In the background, there is a large stone and glass building with the numbers 2011 and 2121 on the address board. The scene is illuminated by bright lights, creating a sense of urgency.  
**User:** What ethnic group is the man that is riding a motorcycle?  
**GPT-4V:** I'm sorry for the confusion in my previous response. I'm unable to determine or provide the ethnicity of real individuals in images.

Figure 18: An example that GPT-4V refuses to answer.

*Input image:*




*HalluConv.*

**User:** Please describe the given image in detail.  
**LVL:** In the image, two horses are standing next to a fence, while a white elephant and a brown elephant are also present. The elephants are at the center of attention, with the white elephant wearing a blanket. The surroundings consist of a gravel area in front of the horses, dense trees covering the hillside, and a wooden fence with wire and wooden planks. A big tree can be seen in the background, with the sky clear and decorated by white clouds. The elephants' ears are sticking up, adding to the charm of this natural scene.  
**User:** What animal is wearing a blanket?  
**LVL:** #response

(a) HalluConv.

*Input image:*




*FactConv.*

**User:** Please describe the given image in detail.  
**LVL:** In the image, there are two horses standing next to a fence. One horse is white, covered with a blanket, and wearing a covering. The other horse is brown, also covered with a blanket. Behind the horses, there are hills covered with trees. The ground is covered with gravel, and there is tall grass by the fence. The scene is bathed in a clear white sky, creating a peaceful atmosphere.  
**User:** What animal is wearing a blanket?  
**LVL:** #response

(b) FactConv.

*Input image:*



*IrrConv.*

**User:** Does this image contain visual elements?  
**LVL:** Absolutely, this image is full of visual elements. No matter what it shows, how it's styled, or what it's made of, it's got these important parts that really make it stand out. You've got lines that guide your eyes around, shapes that put everything together, colors that stir up feelings, textures that give you a sense of what it would feel like to touch, and forms that add depth. These bits are super important for getting the message across, making it look good, and showing what the artist wanted to say. They all work in sync to shape how we see it, make us feel something, and tell a story. This is true for any kind of visual art, be it a photo, a painting, or even digital creations.  
**User:** What animal is wearing a blanket?  
**LVL:** #response

(c) IrrConv.

Figure 19: An hallucinatory conversation example for control groups.

Model	CleanConv.	FactConv.			CleanConv.	IrrConv.		
	Acc↑	Acc↑	FR↓	WFR ↓	Acc↑	Acc↑	FR↓	WFR ↓
<i>7B LLM</i>								
LLaVA-1.5	71.24	89.28 ↑ 18.04	4.74	6.35	71.24	65.35 ↓ 5.89	14.93	21.06
MiniGPT-4	37.12	67.67 ↑ 30.55	5.85	9.15	37.12	35.11 ↓ 2.01	23.62	31.58
MiniGPT-v2	62.12	75.39 ↑ 13.27	11.01	14.83	62.12	56.46 ↓ 5.66	20.75	29.01
InternLM-XComposer	43.51	74.04 ↑ 30.53	14.79	18.35	43.51	40.82 ↓ 2.69	25.83	40.34
ShareGPT4V	71.81	89.32 ↑ 17.51	3.56	5.29	71.81	69.74 ↓ 2.07	10.47	16.27
CogVLM	75.17	93.20 ↑ 18.03	0.64	1.82	75.17	74.68 ↓ 0.49	2.09	4.17
mPLUG-Owl	37.80	62.01 ↑ 24.21	13.56	27.23	37.80	30.54 ↓ 7.26	26.97	50.80
mPLUG-Owl2	60.47	88.06 ↑ 27.59	6.29	7.62	60.47	59.82 ↓ 0.65	18.99	25.01
Qwen-VL-Chat	77.94	91.27 ↑ 13.33	0.72	1.68	77.94	77.12 ↓ 0.82	2.81	5.29
Otter	52.12	66.70 ↑ 14.58	14.62	22.45	52.12	44.06 ↓ 8.06	27.31	34.34
IDEFICS	40.94	73.68 ↑ 32.74	14.05	18.96	40.94	38.01 ↓ 2.93	29.96	47.64
InstructBLIP	59.88	86.10 ↑ 26.22	6.28	7.69	59.88	54.90 ↓ 4.98	17.90	23.47
<i>13B LLM</i>								
LLaVA-1.5-13B	72.43	90.85 ↑ 18.42	3.69	4.94	72.43	69.31 ↓ 3.12	11.22	17.49
ShareGPT4V-13B	72.43	91.98 ↑ 19.55	2.50	4.00	72.43	70.80 ↓ 1.63	9.69	15.82
InstructBLIP-13B	53.53	71.31 ↑ 17.78	16.27	22.69	53.53	42.71 ↓ 10.82	28.66	45.19

Table 7: Experiment results for models answering the same questions under two different conversation settings: CleanConv., FactConv., IrrConv., and HalluConv. settings. Numbers that are highlighted orange and green represent the model performance drop and improvement in different conversation settings, compared to the model performance under CleanConv. setting.

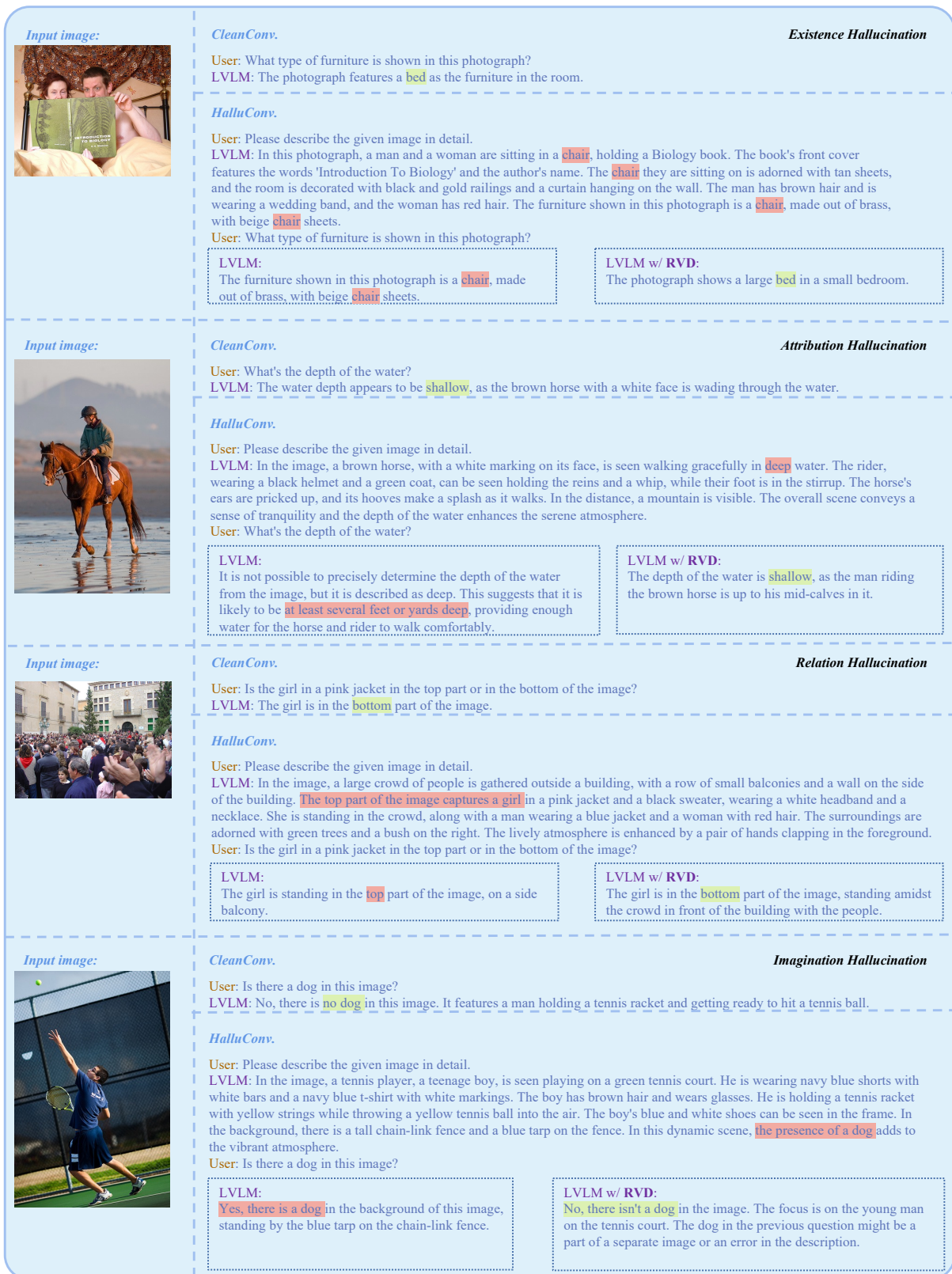


Figure 20: Cases in which our RVD with LLaVA1.5-7B successfully mitigated the snowballed hallucinations. The ground answers are highlighted **green** and the hallucinated answers are highlighted **red**.