

Reefknot: A Comprehensive Benchmark for Relation Hallucination Evaluation, Analysis and Mitigation in Multimodal Large Language Models

Kening Zheng^{1*}, Junkai Chen^{1*}, Yibo Yan^{1,2}, Xin Zou¹, Xuming Hu^{1,2} [†]

¹Hong Kong University of Science and Technology (Guangzhou)

²Hong Kong University of Science and Technology

{neok2zkn, junkai.chen.0917}@gmail.com; {xuminghu}@hkust-gz.com

Abstract

Hallucination issues persistently plagued current multimodal large language models (MLLMs). While existing research primarily focuses on object-level or attribute-level hallucinations, sidelining the more sophisticated relation hallucinations that necessitate advanced reasoning abilities from MLLMs. Besides, recent benchmarks regarding relation hallucinations *lack in-depth evaluation and effective mitigation*. Moreover, their datasets are typically derived from a systematic annotation process, which could introduce *inherent biases* due to the predefined process. To handle the aforementioned challenges, we introduce **Reefknot**, a comprehensive benchmark specifically targeting relation hallucinations, consisting of over 20,000 samples derived from real-world scenarios. Specifically, we first provide a systematic definition of relation hallucinations, integrating perspectives from perceptive and cognitive domains. Furthermore, we construct the relation-based corpus utilizing the representative scene graph dataset Visual Genome (VG), from which semantic triplets follow real-world distributions. Our comparative evaluation across three distinct tasks revealed a substantial shortcoming in the capabilities of current MLLMs to mitigate relation hallucinations. Finally, we advance a novel confidence-based mitigation strategy tailored to tackle the relation hallucinations problem. Across three datasets, including Reefknot, we observed an average reduction of 9.75% in the hallucination rate. We believe our paper sheds valuable insights into achieving trustworthy multimodal intelligence. Our dataset and code will be released upon paper acceptance.

1 Introduction

In recent years, large language models (LLMs) have revolutionized the AI field by expanding their training data to trillions of tokens and increasing their parameter counts to hundreds of billions (Brown et al. 2020; Achiam et al. 2023; Touvron et al. 2023). This has unlocked powerful emergent abilities, and seen widespread applications in diverse domains (Achiam et al. 2023; Yan et al. 2024; Wang et al. 2023a). Recently, the community managed to combine visual backbones with powerful LLMs, resulting in multimodal large language models (MLLMs) (Liu et al. 2023b). While this has led to advancements in multimodal scenarios, it also presents challenges for MLLMs, notably their

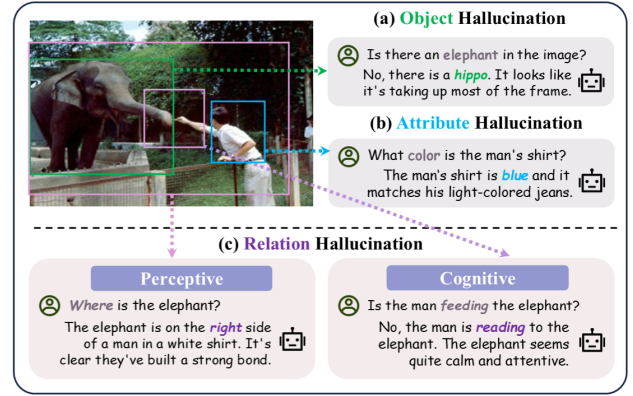


Figure 1: Comparison between the focus of Reefknot benchmark — relation hallucination with two categories (i.e., perceptive and cognitive) vs. object and attribute hallucinations.

tendency to generate hallucinations. In LLMs, hallucinations occur when the model generates inaccurate or misleading factual information that cannot be supported by existing knowledge (Zhang et al. 2023). However, such issues become more complex in MLLMs, as hallucinations can manifest as responses containing references or descriptions of the input image that are incorrect (Bai et al. 2024; Huo et al. 2024a). Therefore, it is crucial to evaluate and mitigate these hallucinations to improve the trustworthiness of MLLMs in real-world scenarios.

Hallucination in MLLMs likely originates from knowledge biases between pre-training and fine-tuning, including statistical biases in the training data, over-reliance on parametric knowledge, and skewed representation learning, as suggested by previous research (Bai et al. 2024; Zhu et al. 2024). Specifically, the hallucination in MLLMs can be divided into three categories: *object*, *attribute* and *relation* hallucinations (Bai et al. 2024; Liu et al. 2024).

As depicted in Figure 1 (a), object-level hallucination focuses on the model's discrimination of the existence of basic objects (Li et al. 2023c); while as shown in Figure 1 (b), attribute-level hallucination often focuses on whether the model can distinguish some properties of the object itself like color, number, shape and so on (Fu et al. 2024). Their

*These authors contributed equally.

[†]Corresponding author.

Benchmarks	Dataset		Evaluation			Analysis	
	Source	Construction	Y/N	MCQ	VQA	Metric	Improv. Focus
POPE (Li et al. 2023c)	COCO	Post-processed	✓	✗	✗	Acc.	Co-occur.
HaELM (Wang et al. 2023b)	MS-COCO	Post-processed	✓	✗	✗	Acc.	Attention
MME (Fu et al. 2024)	Self-Sourced	Manual	✓	✗	✗	Acc.	-
AMBER (Wang et al. 2024a)	MS-COCO	Post-processed	✓	✗	✗	Acc.	-
MHalubench (Chen et al. 2024b)	Self-Sourced	Post-processed	✓	✗	✗	Prec.	-
R-Bench (Wu et al. 2024)	COCO	Post-processed	✓	✗	✗	Acc.	Co-occur.
MMRel (Nie et al. 2024)	Visual-Genome	Synthetic	✓	✗	✗	Acc.	-
VALOR-EVAL (Qiu et al. 2024)	GQA	Post-processed	✗	✗	✗	LLM-based	Co-occur.
Reefknot (Ours)	Visual-Genome	Original	✓	✓	✓	R_{score}	Confidence

Table 1: Comparisons of our proposed Reefknot benchmark with other relevant benchmarks.

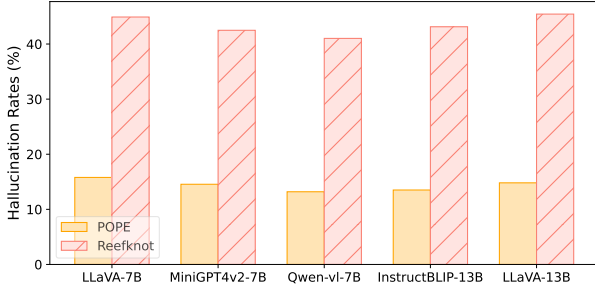


Figure 2: The hallucination rates on POPE (Li et al. 2023c), an object hallucination benchmark, and our Reefknot with a focus on relation hallucination (w/ same configuration).

commonality is that they only focus on the single object present in the image. There has been much work exploring the alleviation of these two types of hallucinations. For example, Woodpecker (Yin et al. 2023) used a post-processing method to correct hallucination after the generation process; VCD (Leng et al. 2023) proposed a contrast decoding strategy to mitigate hallucination by adding noise to images.

Despite these efforts, *the community barely considers relation hallucination, which demands more complex reasoning capabilities from MLLMs*. Figure 1 (c) reveals that relation hallucination is related to at least two objects simultaneously shown in the image, through either perceptive or cognitive perspective. Furthermore, as outlined in Table 1, recent highly relevant benchmarks on hallucinations *lack thorough evaluation and effective mitigation strategies*.

Specifically, these works either adopted either simple Yes/No (Y/N) task to assess the models’ accuracy/precision, or utilized LLM to score the performance, but none of them were able to give a comprehensive evaluation from both discriminative (Y/N MCQ) and generative (VQA) perspectives.

Moreover, previous benchmarks seldom proposed mitigation methods, with only a few focusing on co-occurrence or attention mechanisms to address these issues. In contrast, our paper analyzes token-level confidence at each layer to detect and promptly mitigate hallucinations.

To handle the aforementioned research gaps, we propose the first comprehensive benchmark **Reefknot** to evaluate the

performance on *relation hallucination*. Unlike many benchmarks of MLLM that were constructed by automatic labeling technique, we construct the relation-based dataset based on semantic triplets retrieved from the scene graph dataset. Table 1 also demonstrates that our triplets are from real-life scenarios, without any post-processing (e.g., segmentation and bounding box techniques), manual annotation, and synthetic method (e.g., diffusion-based generation). We categorize relation hallucinations into two types: perceptive, involving concrete relational terms like “on”, “in”, “behind” and cognitive, which encompasses more abstract terms such as “blowing” and “watching”. Second, we evaluate the mainstream MLLMs on Reefknot via three diverse tasks across two types of relation hallucinations. Figure 2 also illustrates that relation hallucination can be more severe than object hallucination in current MLLMs, highlighting the importance of our evaluation. Furthermore, we propose a simple relation hallucination mitigation method named Detect-then-Calibrate. This originates from the experimental observation that when relation hallucinations occur, the response probability drops significantly, hovering just above 50% in extreme cases compared to the usual nearly 90%. Our method achieves an average improvement of 9.75% across three relation hallucination benchmarks.

In summary, Our main contributions are as follows:

1. We have constructed Reefknot, a benchmark comprising two types of relationships and three evaluation tasks to assess relation hallucinations comprehensively.
2. We have conducted a thorough evaluation of relation hallucinations across mainstream MLLMs, uncovering that these models are disproportionately susceptible to perceptual hallucinations in comparison to cognitive ones.
3. We have proposed a novel Detect-then-Calibrate method to detect and mitigate hallucination. By analyzing token confidence scores, we established a threshold to identify hallucinations. Further, we applied a calibration strategy to mitigate hallucination at intermediate confidence levels. Extensive experiments on three relation hallucination datasets demonstrate the effectiveness of our approach.

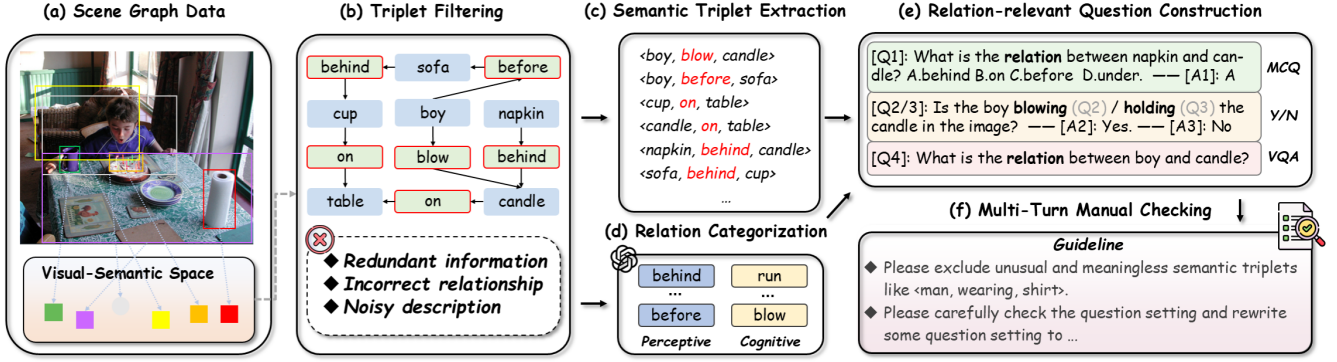


Figure 3: The data construction pipeline of our proposed Reefknot benchmark.

2 Relation Hallucination Benchmark

In this section, we describe the dataset construction pipeline of Reefknot benchmark as shown in Figure 3. Unlike object and attribute hallucinations that only involve one entity, relations involve three entities making it more difficult to handle. We first identify relation triplets from Visual Genome (VG) dataset (Krishna et al. 2016) (Phase a), and conduct triplet filtering (Phase b). Subsequently, we extract the semantic triplets (Phase c) and categorize their relations (Phase d). Then, a relation-based question set can be constructed into three types (Phase e). Finally, the quality of dataset is ensured by three rounds of expert-based validation (Phase f).

Triplet Identification, Filtering and Extraction The dataset comprises of 11,084 images taken from VG dataset, a finely annotated scene graph dataset utilized by the research community (Tang et al. 2020; Liang et al. 2019). As indicated in Figure 3 (a), visual objects and their relations from VG dataset can be easily identified. Besides, we filter the triplets with redundant information, incorrect relationships, or noisy descriptions, as depicted in Figure 3 (b). Subsequently, we can extract semantic triplets by identifying subject-object pairs and the relationships between them, forming (*subject*, *relation*, *object*) triplets in Figure 3 (c).

Relation Categorization As depicted in Figure 3 (d), we categorize relationships into two categories based on deep semantic meanings: perceptive and cognitive. Perceptive relationships involve locational prepositions, such as <boy, behind, sofa> and <cup, on, table>; whereas cognitive relationships are expressed through action phrases indicating states, such as <boy, eating, food> and <girl, sleeping in, bed>. ChatGPT is employed to assist in this classification. Table 2 also indicates the sample numbers in different tasks and hallucination categories. The prompt we used for relation categorization is depicted in Figure 4.

Relation-Relevant Question Construction As shown in Figure 3 (e), we construct three types of relation-relevant question sets to evaluate the state-of-the-art MLLMs’ abilities in relation-level perception and reasoning.

- For Yes/No (Y/N) questions, we employ an adversarial approach by introducing a negative sample within the

Prompt:

You are a relation term classification assistant. Please help me determine whether the following relational terms/phrases belong to perceptive relationships or cognitive relationships. Perceptive relationships are defined as those involving locational or state-based prepositions, such as “on” and “behind”. Cognitive relationships are defined as some words that indicate an action or behavior, such as holding, watching, etc. Here I will give you some demonstrations for reference: {In-context examples} Please answer with a single word in <Perception, Cognition>.

Input: {Input}

Figure 4: Prompt template for relation categorization.

same triplet, alongside the positive sample, to test the model’s ability to correctly answer “No”.

- Multiple Choice Questions (MCQ) are designed with one correct answer and three random options to evaluate the model’s resistance to relation hallucinations within a controlled vocabulary.
- Visual Question Answering (VQA) is an open-ended task that allows us to comprehensively assess a model’s instruction-following capabilities and relation perception within an open-domain environment.

Multi-Turn Manual Checking Finally, we perform multi-turn manual verification to ensure the quality of the question sets (see Figure 3 (f)). Each question undergoes at least three rounds of review by four domain experts. We revise any inappropriate expressions and exclude meaningless questions, such as “Is the window on the wall?”, which lack informative value and can be answered without visual input. After rigorous screening process, our dataset comprises 21,880 questions across 11,084 images as shown in Table 2.

Category	Y/N	MCQ	VQA	Total
#Perception	5,440	4,800	2,150	13,260
#Cognition	4,300	2,150	2,720	8,600
#Total	9,740	6,950	4,870	21,880
Ratio of positive and negative samples				1:1
Number of perceptive relationship types				56
Number of cognitive relationship types				152
Number of images				11084

Table 2: Statistical overview of Reefknot benchmark dataset.

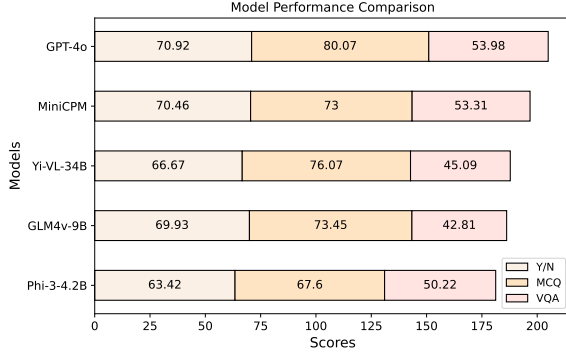


Figure 5: MLLMs with **top five best performance evaluated on Reefknot benchmark**. We report the sum of the respective metric across three tasks for reference.

3 Hallucination Evaluation

3.1 Task and Evaluation Metrics

As illustrated in Table 3, we conducted evaluations on mainstream MLLMs to evaluate relation-level hallucinations. All experiments of open-sourced MLLMs were conducted using 8 NVIDIA A100 GPUs. To reduce variability, each experiment was run three times, and we reported the average of these results. In our evaluation, we reported our results from two distinct categories: perception and cognition.

For discriminative Y/N and MCQ tasks, we reported the hallucination rate $Halr$ as a metric. Generative questions have always posed a challenge in evaluation of MLLMs. In our assessment of the generative VQA task, we employed the DeBERTa model using a bidirectional entailment approach for label match (Kuhn, Gal, and Farquhar 2023), and we denoted $Halr$ as metric for simplicity as well. For more details, we have listed in Appendix. In general, we use the following metric R_{score} to comprehensively evaluate the overall performance across the three tasks:

$$R_{score} = \text{Avg} \left[\sum_{i=1}^3 (1 - Halr_i) \right]. \quad (1)$$

3.2 Main results

Overall Performance Significant performance differences among the various models are evident. Table 3 shows significant differences in the performance of various models across different question types. For instance, Qwen-vl-

chat excels in Y/N and MCQ settings but encounters serious hallucination issues in VQA tasks. A detailed review of the model’s responses reveals that while Qwen can follow instructions accurately, many responses contain expressions that are completely unrelated to the labels. We presume that such a phenomenon is owing to over-fitting training during the instruction-tuning process. Among open-sourced models, MiniCPM (Hu et al. 2024) stands out, likely due to its adoption of the fine-grained alignment technique such as RLHF-V (Yu et al. 2023) to alleviate hallucinations.

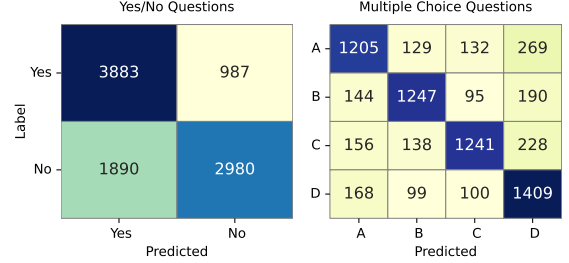


Figure 6: **Confusion matrixes** of MiniCPM-7B on Reefknot benchmark (Left: Y/N setting; Right: MCQ setting).

Error Analysis Figure 5 presents a comparative analysis of the top five best models’ performances on the Reefknot benchmark. The comparison reveals a clear hierarchy in performance, with GPT-4o outperforming the rest across all tasks. MiniCPM, Yi-VL-34B, and GLM4v-9B show a competitive edge, closely following GPT-4o, while Phi-3 lags behind relatively. In terms of task performance, Figure 5 indicates that while some models excel in one task, they may not perform as well in others. For instance, Yi-VL-34B has the second-highest score in MCQ but falls short in VQA. This highlights the importance of task-specific model tuning and the need for a balanced approach to improve overall performance across all relation hallucination tasks.

Perceptive vs. Cognitive Hallucination We find the occurrence of cognitive hallucinations is generally lower than that of perceptive hallucinations, which may diverge from intuition. Across all models and settings, the incidence of perceptive hallucinations is consistently 10% higher than that of cognitive hallucinations. In the most extreme case, such as LLaVA-13B model in the MCQ setting, the rate of perceptive hallucinations is 30.16% higher than that of cognitive hallucinations. A possible reason for this phenomenon is that during the pre-training and fine-tuning processes, models often utilize large-scale image-caption datasets. These datasets typically contain detailed visual descriptions, enabling models to perform better in cognitive relationships such as *running*, *eating*, etc. Conversely, these models may struggle when dealing with some perceptive relationships based on common sense because they were ignored in the annotation process of the original dataset.

To analyze the error cases quantitatively, we visualized the results of MiniCPM-7B, the best open-sourced model, across two discriminative settings, as illustrated in Figure 6.

MLLMs	Size	Perception↓			Cognition↓			Total↑
		Y/N	MCQ	VQA	Y/N	MCQ	VQA	R_{score}
Phi-3-vision-128k-instruct (Abdin et al. 2024)	4.2B	39.88	57.07	50.98	33.97	21.35	49.45	60.30
Yi-VL (AI et al. 2024)	6B	33.56	47.53	71.02	33.16	16.33	74.96	55.81
LLaVA (Liu et al. 2023a)	7B	37.67	68.05	52.93	33.99	51.04	54.56	51.41
MiniGPT4-v2 (Chen et al. 2023)	7B	46.7	78.00	61.30	43.73	68.50	65.88	39.88
MiniCPM (Hu et al. 2024)	7B	31.93	48.65	47.63	27.65	16.71	45.96	65.73
Qwen-VL (Bai et al. 2023)	7B	42.21	56.7	72.47	33.53	21.88	73.01	52.55
Deepseek-VL (DeepSeek-AI et al. 2024)	7B	37.58	56.33	67.07	32.22	23.60	59.34	56.39
GLM4V (GLM et al. 2024)	9B	34.09	50.47	58.09	27.08	16.87	56.47	62.03
LLaVA (Liu et al. 2023a)	13B	40.7	<u>59.35</u>	48.93	34.19	<u>29.19</u>	54.45	57.47
CogVLM (Wang et al. 2024b)	19B	37.23	47.95	70.14	29.89	18.54	66.18	57.1
Yi-VL (AI et al. 2024)	34B	32.79	44.19	57.67	33.75	14.85	52.72	62.61
GPT-4o (OpenAI et al. 2024)	-	32.56	40.93	42.70	26.27	11.53	48.78	68.32

Table 3: **Evaluation of hallucination rates** on the different MLLMs. Additionally, we use **bold** to highlight the best performance of open-sourced MLLMs, and underline to emphasize the distinction between perception and cognition of LLaVA-13B.

For Y/N questions, the model tends to favor positive responses (i.e., Yes). Specifically, among all misclassifications, the instances where a No label is incorrectly classified as Yes are twice as times as instances where a Yes label is incorrectly classified as No. Besides, the model tends to answer D in MCQ settings. We suspect the aforementioned tendency is likely due to the imbalance in the distribution of the training data.

4 Analysis of Relation Hallucination

To quantitatively compare the decision probability distribution when hallucinations occur, we calculate the average probability of an equal number of relation hallucination and non-hallucination examples, as shown in Table 4.

Dataset	Reefknot	MMRel	Rbench
LLaVA	0.67	0.76	0.80
MiniGPT4-v2	0.76	0.75	0.62

Table 4: The **average probability of all hallucination cases**. We test LLaVA and MiniGPT4-v2 on Reefknot and two other representative relation hallucination benchmarks.

The table shows that when hallucinations occur, the confidence level in the answers is quite low. Specifically, experiments conducted on three relational-level datasets indicate that the overall probability of the answers is only about 70%. In contrast, under normal circumstances, large language models can achieve probability values of up to 95% when providing factual and truthful answers. Therefore, a straightforward approach to detecting relation hallucinations is to utilize the entropy $E(X)$ of the probability distribution.

$$E(X) = - \sum_{i=1}^n p(x_i) \log p(x_i). \quad (2)$$

Note that because MLLMs have an extensive vocabulary, it becomes challenging to discern meaningful patterns in entropy variation when predicting the next word across the

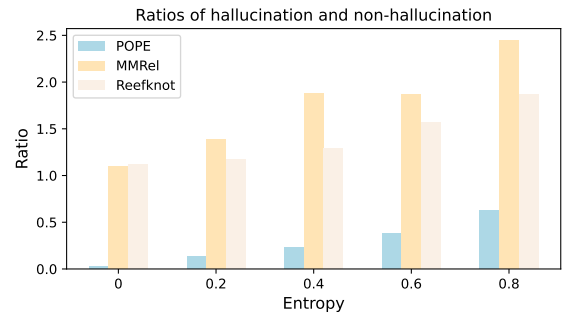


Figure 7: The respective **ratios between hallucination and non-hallucination** with different entropy values. We conducted experiments on POPE, MMRel and our Reefknot.

entire distribution. Consequently, our analysis is restricted to observing the variation patterns of vocabulary within the range of potential answers x_i . We present the ratio of hallucination cases to non-hallucination cases among both object and relation hallucination benchmarks in Figure 7. When $E(X) > 0.6$, *relation hallucinations* occur to a significant degree, indicating the effectiveness of our method to detect *relation hallucination* via entropy. To investigate the mechanism behind hallucination generation, we conducted an in-depth analysis of the fluctuations in probability values across model layers. MLLMs consist of a vision encoder, a projection layer, and a strong LLM decoder, which is stacked by a set of N transformer layers, and an MLP layer $\phi(\cdot)$ for predicting the next-word probability of distribution auto-regressively. Given an image represented by vision encoder as $\mathcal{V}_t = \{v_1, v_2, \dots, v_t\}$ and a text prompt of tokens $\mathcal{P}_t = \{p_1, p_2, \dots, p_t\}$, they are processed as a sequence $\mathcal{H}_0 = \psi(\mathcal{V}_t, \mathcal{P}_t)$ through projection and concatenation function $\psi(\cdot)$. Thus $\mathcal{H}_0 = \{h_1^{(0)}, \dots, h_{t-1}^{(0)}\}$, in which $h_i^{(k)}$ means the hidden states of i_{th} token in k_{th} language decoder layer.

Then \mathcal{H}_0 would be processed by each of the transformer

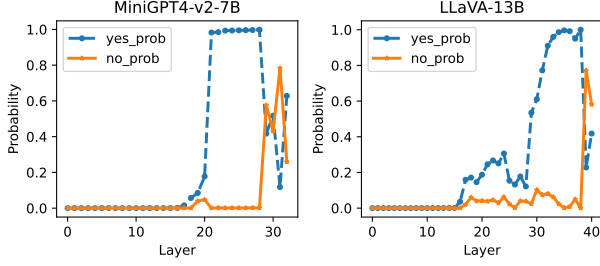


Figure 8: The **average probability among layers when hallucinations occur**. Sharp changes in probability happen in deep layers. Note we report MiniGPTv2-7B with 32 layers and LLaVA-13B with 40 layers to prove universality.

layers in the language decoder successively. We denote the output of the j -th layer as \mathcal{H}_j . In the normal forward process, \mathcal{H}_j will be calculated by \mathcal{N} times, then it will be passed through language model head layer $\phi(\cdot)$ to predict the probability of the next token r_t over the vocabulary set \mathcal{X} . In our set, we manually pass every \mathcal{H}_n to explore the mechanism of hallucination with probability distributions and the next token to generate. For every layer, we can obtain the probability distribution \mathbb{P} and the next word prediction $r_t^{(j)}$.

$$\mathbb{P}(\mathcal{H}_j | \mathcal{H}_{j-1}) = \text{softmax}(\phi(\mathcal{H}_{j-1})), \quad j \in \mathcal{N}. \quad (3)$$

$$r_t^{(j)} = \arg \max \mathbb{P}(\mathcal{H}_j | \mathcal{H}_{j-1}^{(j)}), \quad r_t^{(j)} \in \mathcal{X}. \quad (4)$$

Using Equation 4, we visualized these changes during the forward propagation process, as illustrated in Figure 8. To avoid variability from individual instances, we reported the average values of all data involving hallucinations occurring in Reefknot. Furthermore, to ensure a fair comparison, we employed both the 32-layer MiniGPT4-v2-7B (Chen et al. 2023) and the 40-layer LLaVA-13B (Liu et al. 2023a).

It can be observed that in the shallow layers (0th-20th), the probability of possible answers does not increase. We hypothesize that this is because, in shallow layers, the model is aggregating information to generate an answer. Hallucination occurs in the deep layers. Since deep layers contain a vast amount of knowledge, we speculate that they may cause relation hallucinations.

5 Detect-Then-Calibrate Mitigation Method

The evaluation results indicate that MLLMs commonly suffer from severe *relation hallucinations*. As analyzed in Section 4, we found that these *relation hallucinations* primarily stem from the model’s lack of confidence. The lack of confidence in the model’s responses results in a relatively high entropy value. Therefore, we can detect the occurrence of *relation hallucinations* by monitoring changes in entropy.

Initially, we set a specific entropy threshold γ to detect potential hallucinations in the model’s output. If the entropy of the model’s response exceeds γ , which suggests a significant lack of confidence, we infer a high probability that the model has generated hallucinations. In cases where the model is identified to hallucinate, we will utilize the hidden states

from intermediate layers to calibrate the final outputs layers, which is inspired by DoLa (Chuang et al. 2023) approach. Note, unlike traditional contrastive decoding strategies (Li et al. 2023a; Leng et al. 2023; Huang et al. 2023), we do not calibrate all cases. Instead, we focus on mitigating potential hallucination cases to avoid altering non-hallucinatory cases into hallucination ones. Formula 5 show our calibration process.

$$r = \begin{cases} \arg \max \log \frac{(1+\alpha) \cdot \text{softmax}(\phi(h_t^n))}{\alpha \cdot \text{softmax}(\phi(h_t^{n-\lambda}))}, & \text{if } E_t > \gamma. \\ \arg \max (\text{softmax}(\phi(h_t^n))), & \text{otherwise.} \end{cases} \quad (5)$$

Note that λ is the hyperparameter to control the degree of the intermediate layer. The overall algorithmic flow can be seen from the pseudo-code below:

Algorithm 1: Detect-Then-Calibrate Algorithm

Require: Image v_t ; MLLM $\mathcal{M}(\cdot)$; Prompt p_t ; Uncertainty Entropy threshold γ

- 1: $\mathcal{H}, r_0, \mathbb{P} = \mathcal{M}(x_t, p_t)$
- 2: The entropy of generate token
 $E(r_0) = -\sum_{i=1}^n \rho_i \log \rho_i \quad \rho \in \mathbb{P}$
- 3: **if** $E(r_0) \geq \gamma$ **then**
- 4: Hallucination occurs! Calibrate \mathcal{H} by first λ layer
- 5: $h_\delta = \log \frac{(1+\alpha) \cdot \text{softmax}(\phi(h_t^n))}{\alpha \cdot \text{softmax}(\phi(h_t^{n-\lambda}))}$
- 6: Calibrated response $\bar{r} = \arg \max h_\delta$
- 7: **else**
- 8: Normal response $r = \arg \max (\text{softmax}(\phi(h_t^n)))$
- 9: **end if**

As illustrated in Table 5, we conducted experiments using LLaVA-13B. To demonstrate the robustness of our results, we employed two additional relation hallucination datasets, MMRel (Nie et al. 2024) and R-bench (Wu et al. 2024). During experiments, we set $\lambda = 2, \alpha = 0.1, \gamma = 0.9$. For a fair comparison, in addition to reporting the baseline model, we also report some training-free methods such as VCD (Leng et al. 2023), DoLa (Chuang et al. 2023), and OPERA (Huang et al. 2023). Our approach achieved improvements across all three relation hallucination datasets. Specifically, on the MMRel dataset, our model achieved a 19.7% improvement compared to the baseline setting.

Methods	Reefknot	MMRel	R-bench
Baseline	37.06	40.43	29.52
+ VCD (Leng et al. 2023)	38.32	41.96	22.05
+ DoLa (Chuang et al. 2023)	36.96	39.68	23.52
+ OPERA (Huang et al. 2023)	35.73	39.22	26.73
+ Detect-then-Calibrate (Ours)	34.50	21.73	22.02

Table 5: The **hallucination rates of LLaVA-13B and its variants** with hallucination mitigation techniques across Reefknot and two other relation hallucination benchmarks.

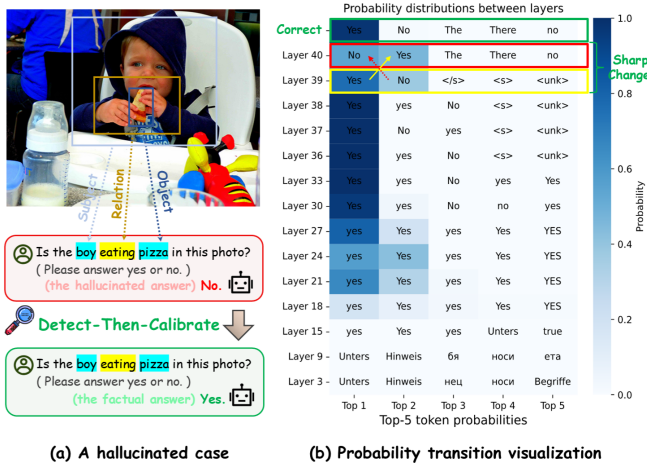


Figure 9: Part (a) illustrates a genuine case of relation hallucination derived from the Reefknot benchmark. The visualization of probability transitions across the layers of the language model is presented in Part (b). Subsequently, our refined results are distinctly highlighted in green, juxtaposed against significant probability variations, which are clearly demarcated within a red-bordered box.

6 Case study

To intuitively demonstrate the generation process of hallucination, we visualized a real case from Reefknot in Figure 9. In the image, when we ask “Is the boy eating pizza in the photo?”, MLLMs are unable to provide precise answers. To investigate this, we analyze the probabilities associated with the top five predicted tokens at each layer of the language model. As illustrated in Figure 9, the model typically converges on the correct answer as it progresses from shallow to deep layers. Notably, hallucinations didn’t occur until the final layer.

Additionally, we can observe that as the number of layers increases, the model increasingly focuses on narrowing down the range of answer choices, with the probability for both ‘yes’ and ‘no’ rising.

However, as depicted in the changes of colors, when it reaches the final Transformer block (layer 38-40th), the model’s choices suddenly become uncertain and entropy rises, accompanied by the emergence of hallucinations. So we used logits of the intermediate layers, which are less likely to hallucinate, to calibrate the final logits to correct the answer. Following our calibration utilizing Equation 5, we not only dispelled any misconceptions but also significantly bolstered the confidence of our responses.

7 Related Work

Hallucination Benchmarks in MLLMs Large Language Models (LLMs) are powerful tools, and exploring their interpretability (Han et al. 2024; Huo et al. 2024b), operational mechanisms, and methods to ensure trustworthy responses (Hu et al. 2023) has become a crucial area of research. In the realm of MLLMs, hallucinations are typically

categorized into three distinct types: object, attribute, and relation, as outlined in prior studies (Yin et al. 2024; Liu et al. 2024). At the object and attribute levels, a considerable number of representative benchmarks such as POPE (Li et al. 2023c) and HaELM (Wang et al. 2023b) have been introduced by researchers. Many benchmarks are accompanied by a diverse array of evaluation criteria, with generative criteria including CHAIR (Li et al. 2023b), THRONE (Kaul et al. 2024) and various discriminative criteria being particularly notable.

The existing relation hallucinations benchmarks (Wu et al. 2024; Nie et al. 2024) focus only on discriminative criteria. The gap of current benchmarks is to identify relationship pairs that accurately reflect the dynamic interactions. Thus, we propose the Reefknot benchmark, incorporating both generative and discriminative criteria to comprehensively evaluate the relation hallucination.

Scene Graph Dataset The concept of scene graphs was first introduced by Zhu et al. (2022), offering a novel method to describe and retrieve complex scenes within images. Building on this foundation, Krishna et al. (2016) further advanced the field by introducing the Visual Genome (VG) dataset in 2017. In recent years, there has been significant progress in leveraging this triplet-based structure for various applications. For instance, Yao et al. (2018) utilized scene graphs to generate more detailed image descriptions.

Confidence Calibration Confidence estimation (Zou et al. 2023; Chen et al. 2024a) and calibration are essential for enhancing the reliability of LLMs such as GPT-3 (Brown et al. 2020). To assess the confidence associated with outputs from LLMs, Kuhn, Gal, and Farquhar (2023) have developed a method called semantic entropy that utilizes linguistic invariances reflecting shared meanings. However, this method relies on accessing token-level probabilities, which are often unavailable through current black-box APIs. Kadavath et al. (2022) have designed prompts that encourage the models to self-assess their responses and to explicitly calculate the probability that an answer is true; while Lin, Hilton, and Evans (2022) have prompted LLMs to provide both an answer and an accompanying confidence level. Additionally, Manakul, Liusie, and Gales (2023) introduced a sampling-based method to identify hallucinated facts in model outputs. In a similar vein, Kadavath et al. (2022) have designed prompts that encourage the models to self-assess their responses and to explicitly calculate the probability that an answer is true, further fine-tuning the model to improve the accuracy of these probabilities. Concurrently, Lin, Hilton, and Evans (2022) have prompted LLMs to provide both an answer and an accompanying confidence level. Manakul, Liusie, and Gales (2023) introduced a sampling-based method to identify hallucinated facts in model outputs.

8 Conclusion and Future Work

In conclusion, we propose a comprehensive benchmark called Reefknot to evaluate and mitigate relation hallucinations in MLLMs. We construct the dataset with over 20k data through a scene graph-based construction pipeline, covering

two discriminative tasks (Y/N and MCQ) and one generative task (VQA). Our in-depth evaluation highlights a substantial performance gap on relation hallucination in existing MLLMs, emphasizing the need for more sophisticated reasoning capabilities. Subsequently, we discover that relation hallucinations tend to occur when MLLMs respond with low confidence. Therefore, we propose a Detect-then-Calibrate method to mitigate the relation hallucination via entropy threshold, with an average reduction of 9.75% in the hallucination rate across Reefknot and two other representative relation hallucination datasets. In general, we anticipate that our proposed Reefknot benchmark will pave the way for future advancements in trustworthy multimodal intelligence.

Despite promising, our proposed approach focuses on mitigating basic discriminative hallucinations, but relation hallucinations in open domains are still challenging to quantitatively assess and mitigate. In future research, we will delve deeper into the underlying causes of hallucinations in open domains and investigate both the mechanisms and mitigation strategies. We anticipate that Reefknot will further improve the reliability and practical utility of MLLMs.

References

- Abdin, M.; Jacobs, S. A.; Awan, A. A.; Aneja, J.; Awadallah, A.; Awadalla, H.; Bach, N.; Bahree, A.; Bakhtiari, A.; Bao, J.; Behl, H.; Benhaim, A.; Bilenko, M.; Bjorck, J.; Bubeck, S.; Cai, Q.; Cai, M.; Mendes, C. C. T.; Chen, W.; Chaudhary, V.; Chen, D.; Chen, D.; Chen, Y.-C.; Chen, Y.-L.; Chopra, P.; Dai, X.; Giorno, A. D.; de Rosa, G.; Dixon, M.; Eldan, R.; Fragoso, V.; Iter, D.; Gao, M.; Gao, M.; Gao, J.; Garg, A.; Goswami, A.; Gunasekar, S.; Haider, E.; Hao, J.; Hewett, R. J.; Huynh, J.; Javaheripi, M.; Jin, X.; Kauffmann, P.; Karampatziakis, N.; Kim, D.; Khademi, M.; Kurilenko, L.; Lee, J. R.; Lee, Y. T.; Li, Y.; Li, Y.; Liang, C.; Liden, L.; Liu, C.; Liu, M.; Liu, W.; Lin, E.; Lin, Z.; Luo, C.; Madan, P.; Mazzola, M.; Mitra, A.; Modi, H.; Nguyen, A.; Norick, B.; Patra, B.; Perez-Becker, D.; Portet, T.; Pryzant, R.; Qin, H.; Radmilac, M.; Rosset, C.; Roy, S.; Ruwase, O.; Saarikivi, O.; Saied, A.; Salim, A.; Santacroce, M.; Shah, S.; Shang, N.; Sharma, H.; Shukla, S.; Song, X.; Tanaka, M.; Tupini, A.; Wang, X.; Wang, L.; Wang, C.; Wang, Y.; Ward, R.; Wang, G.; Witte, P.; Wu, H.; Wyatt, M.; Xiao, B.; Xu, C.; Xu, J.; Xu, W.; Yadav, S.; Yang, F.; Yang, J.; Yang, Z.; Yang, Y.; Yu, D.; Yuan, L.; Zhang, C.; Zhang, C.; Zhang, J.; Zhang, L. L.; Zhang, Y.; Zhang, Y.; Zhang, Y.; and Zhou, X. 2024. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. *arXiv:2404.14219*.
- Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- AI, .; ; Young, A.; Chen, B.; Li, C.; Huang, C.; Zhang, G.; Zhang, G.; Li, H.; Zhu, J.; Chen, J.; Chang, J.; Yu, K.; Liu, P.; Liu, Q.; Yue, S.; Yang, S.; Yang, S.; Yu, T.; Xie, W.; Huang, W.; Hu, X.; Ren, X.; Niu, X.; Nie, P.; Xu, Y.; Liu, Y.; Wang, Y.; Cai, Y.; Gu, Z.; Liu, Z.; and Dai, Z. 2024. Yi: Open Foundation Models by 01.AI. *arXiv:2403.04652*.
- Bai, J.; Bai, S.; Chu, Y.; Cui, Z.; Dang, K.; Deng, X.; Fan, Y.; Ge, W.; Han, Y.; Huang, F.; Hui, B.; Ji, L.; Li, M.; Lin, J.; Lin, R.; Liu, D.; Liu, G.; Lu, C.; Lu, K.; Ma, J.; Men, R.; Ren, X.; Ren, X.; Tan, C.; Tan, S.; Tu, J.; Wang, P.; Wang, S.; Wang, W.; Wu, S.; Xu, B.; Xu, J.; Yang, A.; Yang, H.; Yang, J.; Yang, S.; Yao, Y.; Yu, B.; Yuan, H.; Yuan, Z.; Zhang, J.; Zhang, X.; Zhang, Y.; Zhang, Z.; Zhou, C.; Zhou, J.; Zhou, X.; and Zhu, T. 2023. Qwen Technical Report. *arXiv:2309.16609*.
- Bai, Z.; Wang, P.; Xiao, T.; He, T.; Han, Z.; Zhang, Z.; and Shou, M. Z. 2024. Hallucination of Multimodal Large Language Models: A Survey. *arXiv preprint arXiv:2404.18930*.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901.
- Chen, J.; Zhu, D.; Shen, X.; Li, X.; Liu, Z.; Zhang, P.; Krishnamoorthi, R.; Chandra, V.; Xiong, Y.; and Elhoseiny, M. 2023. MiniGPT-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*.

- Chen, Q.; Qin, L.; Zhang, J.; Chen, Z.; Xu, X.; and Che, W. 2024a. M³CoT: A Novel Benchmark for Multi-Domain Multi-step Multi-modal Chain-of-Thought. *arXiv:2405.16473*.
- Chen, X.; Wang, C.; Xue, Y.; Zhang, N.; Yang, X.; Li, Q.; Shen, Y.; Liang, L.; Gu, J.; and Chen, H. 2024b. Unified Hallucination Detection for Multimodal Large Language Models. *arXiv:2402.03190*.
- Chuang, Y.-S.; Xie, Y.; Luo, H.; Kim, Y.; Glass, J.; and He, P. 2023. DoLa: Decoding by Contrasting Layers Improves Factuality in Large Language Models. *arXiv preprint arXiv:2309.03883*.
- DeepSeek-AI; Liu, A.; Feng, B.; Wang, B.; Wang, B.; Liu, B.; Zhao, C.; Dengr, C.; Ruan, C.; Dai, D.; Guo, D.; Yang, D.; Chen, D.; Ji, D.; Li, E.; Lin, F.; Luo, F.; Hao, G.; Chen, G.; Li, G.; Zhang, H.; Xu, H.; Yang, H.; Zhang, H.; Ding, H.; Xin, H.; Gao, H.; Li, H.; Qu, H.; Cai, J. L.; Liang, J.; Guo, J.; Ni, J.; Li, J.; Chen, J.; Yuan, J.; Qiu, J.; Song, J.; Dong, K.; Gao, K.; Guan, K.; Wang, L.; Zhang, L.; Xu, L.; Xia, L.; Zhao, L.; Zhang, L.; Li, M.; Wang, M.; Zhang, M.; Zhang, M.; Tang, M.; Li, M.; Tian, N.; Huang, P.; Wang, P.; Zhang, P.; Zhu, Q.; Chen, Q.; Du, Q.; Chen, R. J.; Jin, R. L.; Ge, R.; Pan, R.; Xu, R.; Chen, R.; Li, S. S.; Lu, S.; Zhou, S.; Chen, S.; Wu, S.; Ye, S.; Ma, S.; Wang, S.; Zhou, S.; Yu, S.; Zhou, S.; Zheng, S.; Wang, T.; Pei, T.; Yuan, T.; Sun, T.; Xiao, W. L.; Zeng, W.; An, W.; Liu, W.; Liang, W.; Gao, W.; Zhang, W.; Li, X. Q.; Jin, X.; Wang, X.; Bi, X.; Liu, X.; Wang, X.; Shen, X.; Chen, X.; Chen, X.; Nie, X.; Sun, X.; Wang, X.; Liu, X.; Xie, X.; Yu, X.; Song, X.; Zhou, X.; Yang, X.; Lu, X.; Su, X.; Wu, Y.; Li, Y. K.; Wei, Y. X.; Zhu, Y. X.; Xu, Y.; Huang, Y.; Li, Y.; Zhao, Y.; Sun, Y.; Li, Y.; Wang, Y.; Zheng, Y.; Zhang, Y.; Xiong, Y.; Zhao, Y.; He, Y.; Tang, Y.; Piao, Y.; Dong, Y.; Tan, Y.; Liu, Y.; Wang, Y.; Guo, Y.; Zhu, Y.; Wang, Y.; Zou, Y.; Zha, Y.; Ma, Y.; Yan, Y.; You, Y.; Liu, Y.; Ren, Z. Z.; Ren, Z.; Sha, Z.; Fu, Z.; Huang, Z.; Zhang, Z.; Xie, Z.; Hao, Z.; Shao, Z.; Wen, Z.; Xu, Z.; Zhang, Z.; Li, Z.; Wang, Z.; Gu, Z.; Li, Z.; and Xie, Z. 2024. DeepSeek-V2: A Strong, Economical, and Efficient Mixture-of-Experts Language Model. *arXiv:2405.04434*.
- Fu, C.; Chen, P.; Shen, Y.; Qin, Y.; Zhang, M.; Lin, X.; Yang, J.; Zheng, X.; Li, K.; Sun, X.; Wu, Y.; and Ji, R. 2024. MME: A Comprehensive Evaluation Benchmark for Multimodal Large Language Models. *arXiv:2306.13394*.
- GLM, T.; ; Zeng, A.; Xu, B.; Wang, B.; Zhang, C.; Yin, D.; Zhang, D.; Rojas, D.; Feng, G.; Zhao, H.; Lai, H.; Yu, H.; Wang, H.; Sun, J.; Zhang, J.; Cheng, J.; Gui, J.; Tang, J.; Zhang, J.; Sun, J.; Li, J.; Zhao, L.; Wu, L.; Zhong, L.; Liu, M.; Huang, M.; Zhang, P.; Zheng, Q.; Lu, R.; Duan, S.; Zhang, S.; Cao, S.; Yang, S.; Tam, W. L.; Zhao, W.; Liu, X.; Xia, X.; Zhang, X.; Gu, X.; Lv, X.; Liu, X.; Liu, X.; Yang, X.; Song, X.; Zhang, X.; An, Y.; Xu, Y.; Niu, Y.; Yang, Y.; Li, Y.; Bai, Y.; Dong, Y.; Qi, Z.; Wang, Z.; Yang, Z.; Du, Z.; Hou, Z.; and Wang, Z. 2024. ChatGLM: A Family of Large Language Models from GLM-130B to GLM-4 All Tools. *arXiv:2406.12793*.
- Han, C.; Xu, J.; Li, M.; Fung, Y.; Sun, C.; Jiang, N.; Abdelzaher, T.; and Ji, H. 2024. Word Embeddings Are Steers for Language Models. *arXiv:2305.12798*.
- Hu, S.; Tu, Y.; Han, X.; He, C.; Cui, G.; Long, X.; Zheng, Z.; Fang, Y.; Huang, Y.; Zhao, W.; Zhang, X.; Thai, Z. L.; Zhang, K.; Wang, C.; Yao, Y.; Zhao, C.; Zhou, J.; Cai, J.; Zhai, Z.; Ding, N.; Jia, C.; Zeng, G.; Li, D.; Liu, Z.; and Sun, M. 2024. MiniCPM: Unveiling the Potential of Small Language Models with Scalable Training Strategies. *arXiv:2404.06395*.
- Hu, X.; Chen, J.; Li, X.; Guo, Y.; Wen, L.; Yu, P. S.; and Guo, Z. 2023. Do Large Language Models Know about Facts? *arXiv:2310.05177*.
- Huang, Q.; Dong, X.; zhang, P.; Wang, B.; He, C.; Wang, J.; Lin, D.; Zhang, W.; and Yu, N. 2023. OPERA: Alleviating Hallucination in Multi-Modal Large Language Models via Over-Trust Penalty and Retrospection-Allocation. *arXiv preprint arXiv:2311.17911*.
- Huo, J.; Yan, Y.; Hu, B.; Yue, Y.; and Hu, X. 2024a. MM-Neuron: Discovering Neuron-Level Domain-Specific Interpretation in Multimodal Large Language Model. *arXiv preprint arXiv:2406.11193*.
- Huo, J.; Yan, Y.; Hu, B.; Yue, Y.; and Hu, X. 2024b. MMNeuron: Discovering Neuron-Level Domain-Specific Interpretation in Multimodal Large Language Model. *arXiv:2406.11193*.
- Kadavath, S.; Conerly, T.; Askell, A.; Henighan, T.; Drain, D.; Perez, E.; Schiefer, N.; Hatfield-Dodds, Z.; DasSarma, N.; Tran-Johnson, E.; Johnston, S.; El-Showk, S.; Jones, A.; Elhage, N.; Hume, T.; Chen, A.; Bai, Y.; Bowman, S.; Fort, S.; Ganguli, D.; Hernandez, D.; Jacobson, J.; Kernion, J.; Kravec, S.; Lovitt, L.; Ndousse, K.; Olsson, C.; Ringer, S.; Amodei, D.; Brown, T.; Clark, J.; Joseph, N.; Mann, B.; McCandlish, S.; Olah, C.; and Kaplan, J. 2022. Language Models (Mostly) Know What They Know. *arXiv:2207.05221*.
- Kaul, P.; Li, Z.; Yang, H.; Dukler, Y.; Swaminathan, A.; Taylor, C. J.; and Soatto, S. 2024. THRONE: An Object-based Hallucination Benchmark for the Free-form Generations of Large Vision-Language Models. *arXiv:2405.05256*.
- Krishna, R.; Zhu, Y.; Groth, O.; Johnson, J.; Hata, K.; Kravitz, J.; Chen, S.; Kalantidis, Y.; Li, L.-J.; Shamma, D. A.; Bernstein, M. S.; and Li, F.-F. 2016. Visual Genome: Connecting Language and Vision Using Crowd-sourced Dense Image Annotations. *arXiv:1602.07332*.
- Kuhn, L.; Gal, Y.; and Farquhar, S. 2023. Semantic Uncertainty: Linguistic Invariances for Uncertainty Estimation in Natural Language Generation. *arXiv:2302.09664*.
- Leng, S.; Zhang, H.; Chen, G.; Li, X.; Lu, S.; Miao, C.; and Bing, L. 2023. Mitigating Object Hallucinations in Large Vision-Language Models through Visual Contrastive Decoding. *arXiv:2311.16922*.
- Li, X. L.; Holtzman, A.; Fried, D.; Liang, P.; Eisner, J.; Hashimoto, T.; Zettlemoyer, L.; and Lewis, M. 2023a. Contrastive Decoding: Open-ended Text Generation as Optimization. In Rogers, A.; Boyd-Graber, J.; and Okazaki, N., eds., *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 12286–12312. Toronto, Canada: Association for Computational Linguistics.

- Li, Y.; Du, Y.; Zhou, K.; Wang, J.; Zhao, W. X.; and Wen, J.-R. 2023b. Evaluating Object Hallucination in Large Vision-Language Models. *arXiv:2305.10355*.
- Li, Y.; Du, Y.; Zhou, K.; Wang, J.; Zhao, X.; and Wen, J.-R. 2023c. Evaluating Object Hallucination in Large Vision-Language Models. In Bouamor, H.; Pino, J.; and Bali, K., eds., *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 292–305. Singapore: Association for Computational Linguistics.
- Liang, Y.; Bai, Y.; Zhang, W.; Qian, X.; Zhu, L.; and Mei, T. 2019. VrR-VG: Refocusing Visually-Relevant Relationships. *arXiv:1902.00313*.
- Lin, S.; Hilton, J.; and Evans, O. 2022. Teaching Models to Express Their Uncertainty in Words. *arXiv:2205.14334*.
- Liu, H.; Li, C.; Li, Y.; and Lee, Y. J. 2023a. Improved Baselines with Visual Instruction Tuning. *arXiv preprint arXiv:2310.03744*.
- Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2023b. Visual Instruction Tuning.
- Liu, H.; Xue, W.; Chen, Y.; Chen, D.; Zhao, X.; Wang, K.; Hou, L.; Li, R.; and Peng, W. 2024. A Survey on Hallucination in Large Vision-Language Models. *arXiv:2402.00253*.
- Manakul, P.; Liusie, A.; and Gales, M. J. F. 2023. SelfCheckGPT: Zero-Resource Black-Box Hallucination Detection for Generative Large Language Models. *arXiv:2303.08896*.
- Nie, J.; Zhang, G.; An, W.; Tan, Y.-P.; Kot, A. C.; and Lu, S. 2024. MMRel: A Relation Understanding Dataset and Benchmark in the MLLM Era. *arXiv preprint arXiv:2406.09121*.
- OpenAI; Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altenschmidt, J.; Altman, S.; Anadkat, S.; Avila, R.; Babuschkin, I.; Balaji, S.; Balcom, V.; Baltescu, P.; Bao, H.; Bavarian, M.; Belgum, J.; Bello, I.; Berdine, J.; Bernadett-Shapiro, G.; Berner, C.; Bogdonoff, L.; Boiko, O.; Boyd, M.; Brakman, A.-L.; Brockman, G.; Brooks, T.; Brundage, M.; Button, K.; Cai, T.; Campbell, R.; Cann, A.; Carey, B.; Carlson, C.; Carmichael, R.; Chan, B.; Chang, C.; Chantzis, F.; Chen, D.; Chen, S.; Chen, R.; Chen, J.; Chen, M.; Chess, B.; Cho, C.; Chu, C.; Chung, H. W.; Cummings, D.; Currier, J.; Dai, Y.; Decareaux, C.; Degry, T.; Deutsch, N.; Deville, D.; Dhar, A.; Dohan, D.; Dowling, S.; Dunning, S.; Ecoffet, A.; Eleti, A.; Eloundou, T.; Farhi, D.; Fedus, L.; Felix, N.; Fishman, S. P.; Forte, J.; Fulford, I.; Gao, L.; Georges, E.; Gibson, C.; Goel, V.; Gogineni, T.; Goh, G.; Gontijo-Lopes, R.; Gordon, J.; Grafstein, M.; Gray, S.; Greene, R.; Gross, J.; Gu, S. S.; Guo, Y.; Hallacy, C.; Han, J.; Harris, J.; He, Y.; Heaton, M.; Heidecke, J.; Hesse, C.; Hickey, A.; Hickey, W.; Hoeschele, P.; Houghton, B.; Hsu, K.; Hu, S.; Hu, X.; Huizinga, J.; Jain, S.; Jain, S.; Jiang, J.; Jiang, A.; Jiang, R.; Jin, H.; Jin, D.; Jomoto, S.; Jonn, B.; Jun, H.; Kafkhan, T.; Łukasz Kaiser; Kamali, A.; Kanitscheider, I.; Kesar, N. S.; Khan, T.; Kilpatrick, L.; Kim, J. W.; Kim, C.; Kim, Y.; Kirchner, J. H.; Kiros, J.; Knight, M.; Kokotajlo, D.; Łukasz Kondraciuk; Kondrich, A.; Konstantinidis, A.; Kosic, K.; Krueger, G.; Kuo, V.; Lampe, M.; Lan, I.; Lee, T.; Leike,
- J.; Leung, J.; Levy, D.; Li, C. M.; Lim, R.; Lin, M.; Lin, S.; Litwin, M.; Lopez, T.; Lowe, R.; Lue, P.; Makanju, A.; Malfacini, K.; Manning, S.; Markov, T.; Markovski, Y.; Martin, B.; Mayer, K.; Mayne, A.; McGrew, B.; McKinney, S. M.; McLeavey, C.; McMillan, P.; McNeil, J.; Medina, D.; Mehta, A.; Menick, J.; Metz, L.; Mishchenko, A.; Mishkin, P.; Monaco, V.; Morikawa, E.; Mossing, D.; Mu, T.; Murati, M.; Murk, O.; Mély, D.; Nair, A.; Nakano, R.; Nayak, R.; Nee-lakantan, A.; Ngo, R.; Noh, H.; Ouyang, L.; O’Keefe, C.; Pachocki, J.; Paino, A.; Palermo, J.; Pantuliano, A.; Parascandolo, G.; Parish, J.; Parparita, E.; Passos, A.; Pavlov, M.; Peng, A.; Perelman, A.; de Avila Belbute Peres, F.; Petrov, M.; de Oliveira Pinto, H. P.; Michael; Pokorny; Pokrass, M.; Pong, V. H.; Powell, T.; Power, A.; Power, B.; Proehl, E.; Puri, R.; Radford, A.; Rae, J.; Ramesh, A.; Raymond, C.; Real, F.; Rimbach, K.; Ross, C.; Rotsted, B.; Roussez, H.; Ryder, N.; Saltarelli, M.; Sanders, T.; Santurkar, S.; Sastry, G.; Schmidt, H.; Schnurr, D.; Schulman, J.; Sel-sam, D.; Sheppard, K.; Sherbakov, T.; Shieh, J.; Shoker, S.; Shyam, P.; Sidor, S.; Sigler, E.; Simens, M.; Sitkin, J.; Slama, K.; Sohl, I.; Sokolowsky, B.; Song, Y.; Staudacher, N.; Such, F. P.; Summers, N.; Sutskever, I.; Tang, J.; Tezak, N.; Thompson, M. B.; Tillet, P.; Tootoonchian, A.; Tseng, E.; Tuggle, P.; Turley, N.; Tworek, J.; Uribe, J. F. C.; Val-lone, A.; Vijayvergiya, A.; Voss, C.; Wainwright, C.; Wang, J. J.; Wang, A.; Wang, B.; Ward, J.; Wei, J.; Weinmann, C.; Welihinda, A.; Welinder, P.; Weng, J.; Weng, L.; Wiethoff, M.; Willner, D.; Winter, C.; Wolrich, S.; Wong, H.; Work-man, L.; Wu, S.; Wu, J.; Wu, M.; Xiao, K.; Xu, T.; Yoo, S.; Yu, K.; Yuan, Q.; Zaremba, W.; Zellers, R.; Zhang, C.; Zhang, M.; Zhao, S.; Zheng, T.; Zhuang, J.; Zhuk, W.; and Zoph, B. 2024. GPT-4 Technical Report. *arXiv:2303.08774*.
- Qiu, H.; Hu, W.; Dou, Z.-Y.; and Peng, N. 2024. VALOR-EVAL: Holistic Coverage and Faithfulness Evaluation of Large Vision-Language Models. *arXiv preprint arXiv:2404.13874*.
- Tang, K.; Niu, Y.; Huang, J.; Shi, J.; and Zhang, H. 2020. Unbiased scene graph generation from biased training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 3716–3725.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Wang, C.; Liu, X.; Yue, Y.; Tang, X.; Zhang, T.; Jiayang, C.; Yao, Y.; Gao, W.; Hu, X.; Qi, Z.; Wang, Y.; Yang, L.; Wang, J.; Xie, X.; Zhang, Z.; and Zhang, Y. 2023a. Survey on Factuality in Large Language Models: Knowledge, Retrieval and Domain-Specificity. *arXiv:2310.07521*.
- Wang, J.; Wang, Y.; Xu, G.; Zhang, J.; Gu, Y.; Jia, H.; Wang, J.; Xu, H.; Yan, M.; Zhang, J.; and Sang, J. 2024a. AMBER: An LLM-free Multi-dimensional Benchmark for MLLMs Hallucination Evaluation. *arXiv:2311.07397*.
- Wang, J.; Zhou, Y.; Xu, G.; Shi, P.; Zhao, C.; Xu, H.; Ye, Q.; Yan, M.; Zhang, J.; Zhu, J.; Sang, J.; and Tang, H. 2023b. Evaluation and Analysis of Hallucination in Large Vision-Language Models. *arXiv:2308.15126*.

Wang, W.; Lv, Q.; Yu, W.; Hong, W.; Qi, J.; Wang, Y.; Ji, J.; Yang, Z.; Zhao, L.; Song, X.; Xu, J.; Xu, B.; Li, J.; Dong, Y.; Ding, M.; and Tang, J. 2024b. CogVLM: Visual Expert for Pretrained Language Models. *arXiv:2311.03079*.

Wu, M.; Ji, J.; Huang, O.; Li, J.; Wu, Y.; Sun, X.; and Ji, R. 2024. Evaluating and Analyzing Relationship Hallucinations in Large Vision-Language Models. *arXiv:2406.16449*.

Yan, Y.; Wen, H.; Zhong, S.; Chen, W.; Chen, H.; Wen, Q.; Zimmermann, R.; and Liang, Y. 2024. Urbanclip: Learning text-enhanced urban region profiling with contrastive language-image pretraining from the web. In *Proceedings of the ACM on Web Conference 2024*, 4006–4017.

Yao, T.; Pan, Y.; Li, Y.; and Mei, T. 2018. Exploring Visual Relationship for Image Captioning. *arXiv:1809.07041*.

Yin, S.; Fu, C.; Zhao, S.; Li, K.; Sun, X.; Xu, T.; and Chen, E. 2024. A Survey on Multimodal Large Language Models. *arXiv:2306.13549*.

Yin, S.; Fu, C.; Zhao, S.; Xu, T.; Wang, H.; Sui, D.; Shen, Y.; Li, K.; Sun, X.; and Chen, E. 2023. Woodpecker: Hallucination Correction for Multimodal Large Language Models. *arXiv:2310.16045*.

Yu, T.; Yao, Y.; Zhang, H.; He, T.; Han, Y.; Cui, G.; Hu, J.; Liu, Z.; Zheng, H.-T.; Sun, M.; et al. 2023. Rlhv-v: Towards trustworthy mllms via behavior alignment from fine-grained correctional human feedback. *arXiv preprint arXiv:2312.00849*.

Zhang, Y.; Li, Y.; Cui, L.; Cai, D.; Liu, L.; Fu, T.; Huang, X.; Zhao, E.; Zhang, Y.; Chen, Y.; Wang, L.; Luu, A. T.; Bi, W.; Shi, F.; and Shi, S. 2023. Siren’s Song in the AI Ocean: A Survey on Hallucination in Large Language Models. *arXiv:2309.01219*.

Zhu, G.; Zhang, L.; Jiang, Y.; Dang, Y.; Hou, H.; Shen, P.; Feng, M.; Zhao, X.; Miao, Q.; Shah, S. A. A.; and Benamoun, M. 2022. Scene Graph Generation: A Comprehensive Survey. *arXiv:2201.00443*.

Zhu, J.; Liu, S.; Yu, Y.; Tang, B.; Yan, Y.; Li, Z.; Xiong, F.; Xu, T.; and Blaschko, M. B. 2024. FastMem: Fast Memorization of Prompt Improves Context Awareness of Large Language Models. *arXiv preprint arXiv:2406.16069*.

Zou, X.; Tang, C.; Zheng, X.; Li, Z.; He, X.; An, S.; and Liu, X. 2023. DPNET: Dynamic Poly-attention Network for Trustworthy Multi-modal Classification. In *Proceedings of the 31st ACM International Conference on Multimedia*, 3550–3559.

We listed more error analyses for Y/N and MCQ for MLLMs. In Figure 11, it shows Yi-VL-34b-chat. Figure 12 shows GLM4v-9b-chat. Figure 13 shows Qwen-vl-chat. Figure 14 shows Phi-vision-128k. It can be observed that, with the exception of Phi-vision-128k, all models demonstrate identical distribution and preference trends.

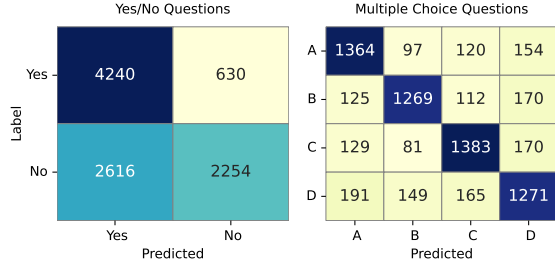


Figure 11: Confusion matrixes of Yi-vl-34b-chat on Reef-knot benchmark (Left: Y/N setting; Right: MCQ setting).

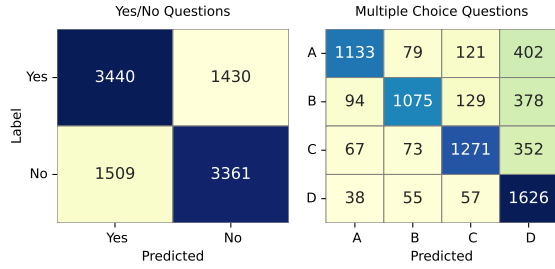


Figure 12: Confusion matrixes of GLM4v-9b-chat on Reef-knot benchmark (Left: Y/N setting; Right: MCQ setting).

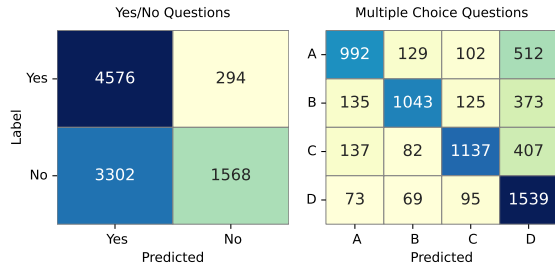


Figure 13: Confusion matrixes of Qwen-vl-chat on Reefknot benchmark (Left: Y/N setting; Right: MCQ setting).

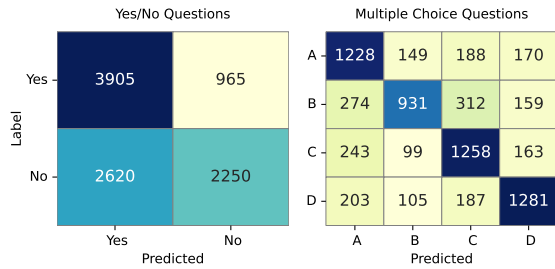


Figure 14: Confusion matrixes of Phi3-vision-128k on Reef-knot benchmark (Left: Y/N setting; Right: MCQ setting).

Visualization of Relation Word

In Figure 18, we present a word cloud that visualizes the proportion of relational terms within our dataset. It can be observed that, due to the use of semantic triples, our relational terms exhibit greater diversity.

VQA Criterion

Here are the evaluation criteria for VQA questions. We use the Deberta model to determine whether the models entail each other. Only when the label and response contain each other will it be judged as a correct reply. We show our key function code in Figure 15. Besides, we show a VQA case in Figure 17.

More Cases

Figure 16 presents examples of Yes/No, multiple-choice questions (MCQ) and visual question answering (VQA) tasks under perceptive and cognitive conditions, shown separately for comparison.

```

1 def are_equivalent(label,response , model, tokenizer, device):
2     #label is groundtruth;response is the response of LLM;
3     def check_implication(label, response):
4         inputs = tokenizer(label,response, return_tensors="pt").to(device)
5         outputs = model(**inputs)
6         logits = outputs.logits
7         largest_index = torch.argmax(F.softmax(logits, dim=1))
8         return largest_index.cpu().item()
9     implication_1 = check_implication(label, response)
10    implication_2 = check_implication(response, label)
11
12
13    assert (implication_1 in [0, 1, 2]) and (implication_2 in [0, 1, 2])
14    implications = [implication_1, implication_2]
15    semantically_equivalent = (implications[0] == 2) and (implications[1] == 2)
16    # only when both are 2, they are semantically equivalent
17    return "yes" if semantically_equivalent else "no"

```

Figure 15: Function to check the semantic equivalence of response and our label.








	Y/N	MCQ	VQA
Perception	 <p>Question: Is the book under laptop in this photo?</p> <p>Answer: Yes</p>	 <p>Question: What is the relation with pizza and tray in this photo?</p> <p>A. on B. besides C. below D. on left of</p> <p>Answer: A</p>	 <p>Question: What is the relation with cat and box in this photo?</p> <p>following format: cat is <relation> box.</p> <p>Answer: Cat is in box</p>
Cognition	 <p>Question: Are the men feeding elep hant in this photo?</p> <p>Answer: No</p>	 <p>Question: What is the relation with cat and plate in this photo?</p> <p>A. cooking B. clutching C. licking D. trashing</p> <p>Answer: C</p>	 <p>Q: What is the relation with woman and zebra in this photo?</p> <p>following format: woman is <relation> zebra.</p> <p>Answer: feeding zebra</p>

Figure 16: Real-world cases from our proposed Reefknof benchmark. We outline the questions for three types of tasks between perception and cognition perspectives as a reference.



What is the relation with cat and laptop in this photo? Please answer in the following format: cat is <relation> laptop

GT Answer: Cat is on laptop.

Cat is sitting on laptop. Cat is under laptop.

Deberta For Evaluate

Correct! Cat is sitting on laptop

Wrong! Cat is not under laptop

Figure 17: An Example of using deberta for evaluation

