

# Investigating and Mitigating Object Hallucinations in Pretrained Vision-Language (CLIP) Models

Yufang Liu<sup>1\*</sup>, Tao Ji<sup>2,3\*</sup>, Changzhi Sun<sup>1</sup>, Yuanbin Wu<sup>1</sup>, Aimin Zhou<sup>1</sup>

<sup>1</sup>School of Computer Science and Technology, East China Normal University

<sup>2</sup>School of Computer Science, Fudan University

<sup>3</sup>Pazhou Laboratory, Huangpu

yfliu.antnlp@gmail.com taoji@fudan.edu.cn ybwu@cs.ecnu.edu.cn

## Abstract

Large Vision-Language Models (LVLMs) have achieved impressive performance, yet research has pointed out a serious issue with object hallucinations within these models. However, there is no clear conclusion as to which part of the model these hallucinations originate from. In this paper, we present an in-depth investigation into the object hallucination problem specifically within the CLIP model, which serves as the backbone for many state-of-the-art vision-language systems. We unveil that even in isolation, the CLIP model is prone to object hallucinations, suggesting that the hallucination problem is not solely due to the interaction between vision and language modalities. To address this, we propose a counterfactual data augmentation method by creating negative samples with a variety of hallucination issues. We demonstrate that our method can effectively mitigate object hallucinations for the CLIP model, and we show that the enhanced model can be employed as a visual encoder, effectively alleviating the object hallucination issue in LVLMs.<sup>1</sup>

## 1 Introduction

Current Large Vision-Language Models (LVLMs) demonstrate significant potential in tasks requiring joint visual and linguistic perception, such as image captioning (Agrawal et al., 2019b), visual question answering (Antol et al., 2015), visual grounding (Yu et al., 2016), and autonomous agents (Durante et al., 2024; Xi et al., 2023). Despite the success of LVLMs, previous studies have revealed that they commonly suffer from hallucinations in practice, including object hallucinations (Li et al., 2023c; Leng et al., 2023; Zhou et al., 2023), spatial hallucinations (Kamath et al., 2023), attribute hallucinations (Zhang et al., 2024), etc. It is widely believed that hallucinations degrade model performance and

reliability, and severely impair the user experience in real-world applications (Ji et al., 2023).

In this work, we focus on investigating the causes of the highly-concerned *object hallucinations*, i.e., LVLMs generate nonexistent objects in the image (Biten et al., 2022). A typical LVLM utilizes a Large Language Model (LLM) as its cognitive foundational model and employs a pre-trained image encoder as its visual perception module (mainly the CLIP encoder). Kamath et al. (2023) investigated the spatial hallucination (e.g., confusing “left of” and “right of”) in LVLMs, and they found that various CLIP encoders struggle to recognize simple spatial relationships (achieving only a 55.0% accuracy on benchmarks, whereas humans are 98.8%). Inspired by their findings, we hypothesize that the CLIP visual encoder might also be one of the causes of object hallucinations.

Hence, we first curate the **Object Hallucination Detection (OHD-Caps)** benchmark from subsets of the COCO (Lin et al., 2014), Flickr30K (Young et al., 2014), and Nocaps (as an out-of-domain benchmark because it comprises unseen objects) (Agrawal et al., 2019a) image caption datasets respectively, to more strictly measure the extent of object hallucinations present in CLIP encoders. We randomly select 16k/1k/1.5k (train/dev/test) samples, with each sample containing one image, one positive descriptive text, and 27 negative descriptive texts. The negative samples are perturbations of the positive sample, achieved by *adding* descriptions of nonexistent objects or *reducing* descriptions of existing objects. Theoretically, a CLIP model without object hallucinations should accurately assign the highest CLIP score to the positive sample. However, taking the most commonly used “CLIP ViT-L/14” in LVLMs as an example, it only scores the highest for positive samples in 19.0% of cases. Since we have observed that the CLIP encoder already has a serious issue with object hallucination, how can we mitigate it?

\* Equal contribution.

<sup>1</sup>Our benchmark and code are publicly available on [https://github.com/Yufang-Liu/clip\\_hallucination](https://github.com/Yufang-Liu/clip_hallucination).

In the contrastive pretraining of CLIP, negative samples come from text descriptions of other images within the batch, which makes the distinction between them quite straightforward. However, mitigating object hallucinations requires the CLIP encoder to be able to differentiate between subtle errors at the object level. We further fine-tune the CLIP model using the training set from **OHD-Caps**. By incorporating a fine-grained object-level contrastive loss, we greatly reduce object hallucinations in the CLIP. Then employing the fine-tuned CLIP as the visual encoder, the object hallucinations in our retrained LVLM, LLaVA-1.5, are also diminished.

In this paper, we study the object hallucinations of CLIP models. Our main contributions are,

- we propose a benchmark, **OHD-Caps**, for evaluating object hallucinations in CLIP models.
- we quantitatively evaluate a wide range of encoders from the CLIP family and find that they all exhibit severe object hallucination issues.
- we propose a fine-grained object-level contrastive loss to further fine-tune the CLIP model, significantly alleviating its object hallucination issues (e.g., from 14.3 to 82.5 for “CLIP ViT-B/32”) and concurrently reducing the hallucination problems of the LLaVA-1.5 (from 80.2 to 83.2 on Nocaps), which uses it as a visual encoder.

## 2 Related Work

### 2.1 Large Vision-Language Model

Recently, inspired by the success of large language models (LLMs), researchers have begun to dedicate efforts to enhance vision language models (VLMs) by integrating robust LLMs, aiming to broaden the knowledge scope of the model and amplify its linguistic comprehension capabilities.

LVLM architectures typically consist of three components: a visual encoder, a modality connection module, and a LLM. The visual encoder and LLM are typically fixed large pretrained models, the visual encoder is usually a variant of the CLIP model (Radford et al., 2021), used for extract visual features, while the LLM, such as LLaMA (Touvron et al., 2023) and Vicuna (Chiang et al., 2023), is used to integrate image information and text information, and completes the prediction of the target. Research focuses on optimizing modality connection modules, with approaches like

Flamingo’s (Alayrac et al., 2022) cross-attention module, LLaVA’s (Liu et al., 2023c) linear layer, and BLIP2’s (Li et al., 2023a) Q-former, diverse yet all boosting VLM performance on various vision-language tasks.

### 2.2 Hallucination in LVLMs

Despite the fact that LVLMs perform well in solving visual-language tasks, they are also plagued by hallucinations. The problem of hallucinations in LVLMs mainly refers to the mismatch between visual input and textual output. For example, in the image captioning task, hallucination refers to the generation of captions that describe objects that do not exist in the image. Although the hallucination problem of LLMs has been widely studied in the NLP field (Ji et al., 2023), there has not been enough research on mitigating the hallucination issue in LVLMs (Shekhar et al., 2017; Liu et al., 2024, 2023a). Recent efforts to mitigate hallucination in LVLMs have focused on enhancing each component of the model. For example, Liu et al. (2023b); Hu et al. (2023) construct instruction-tuning datasets with contrastive question-answer pairs for LVLMs; Sun et al. (2023b); Yu et al. (2023) employ Reinforcement Learning from Human Feedback (RLHF) (Stiennon et al., 2020) to enhance the connection module between the modalities; Leng et al. (2023) propose a visual contrastive decoding strategy for LLM decoing. Despite the wide application of the CLIP model in VLMs and its in-depth study in pairwise comparison context (Yüksekgönül et al., 2023; Hsieh et al., 2023), there has been little discussion on its evaluation regarding hallucinations. Our research addresses this gap in the literature.

## 3 The OHD-Caps Benchmark

Recent studies have found that LVLMs are prone to object hallucinations (Li et al., 2023c; Zhou et al., 2023). In response, researchers have developed several datasets to assess the extent of these hallucinations in such models (Li et al., 2023c; Wang et al., 2023c). However, there is a relative lack of assessment work regarding the hallucinatory effects of the CLIP model, which is widely used as a visual encoder within LVLMs. In this section, we introduce the **Object Hallucination Detection** benchmark (OHD-Caps) we create to evaluate the object hallucination problem in CLIP models and the pipeline for evaluations. Figure 1 shows the

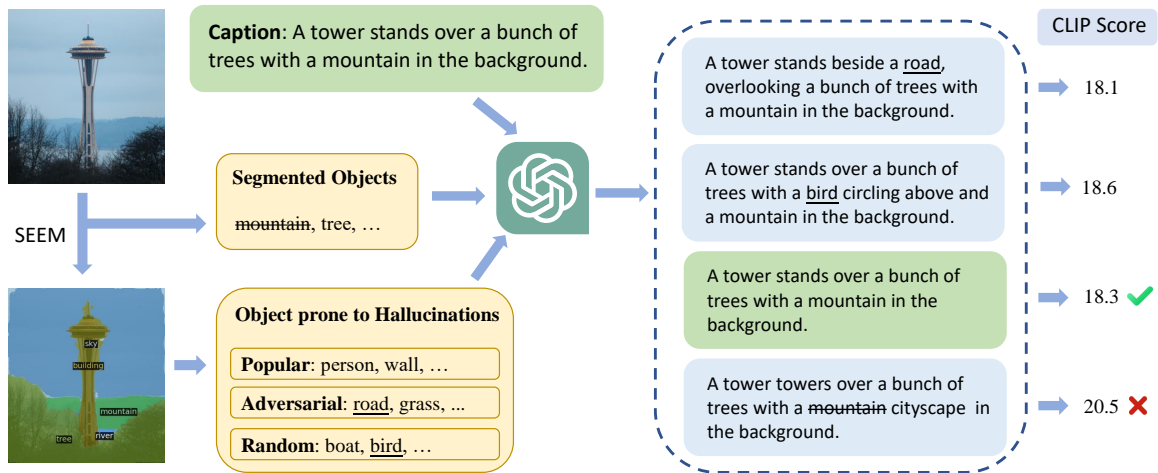


Figure 1: The pipeline of our benchmark creation process. For an image, we first use SEEM (Zou et al., 2023) to identify objects within the image and obtain illusory objects that do not exist in the picture through different sampling strategies. Then we ask GPT to insert or delete objects in the original sentences to create negative samples. We provide both positive and negative samples to the CLIP model to observe if the model predicts the positive samples as having the highest score. This image is from the NoCaps dataset, and the model is CLIP ViT-B/32.

pipeline of our benchmark creation process.

### 3.1 Dataset Construction

CLIP is a versatile neural network that excels at image understanding and can predict text for images in a zero-shot manner. To evaluate the CLIP model’s ability to handle object hallucinations in paired comparison scenarios, given an image with a correct caption, we create incorrect captions containing hallucinatory content. The purpose is to observe whether the model can accurately select the correct text without hallucinations.

**Inserting Hallucinatory Objects** Previous work (Li et al., 2023c; Zhou et al., 2023) show that LVLMs are more prone to generate hallucinatory responses for objects that frequently appear in the dataset. Inspired by this, we create negative samples by inserting objects prone to hallucination into the correct captions. To collect object annotations, we first use SEEM (Zou et al., 2023) to automatically segment objects in the images. Three kinds of hallucinatory objects are collected: *random objects* which are sampled randomly, *popular objects* which are the top frequent objects in the whole dataset, and *adversarial objects* which are the top frequent objects with the segmented objects. Each category contains three objects. To create examples with varying levels of hallucinations, we attempt to insert one to three objects for each category, resulting in each type of hallucination containing a total of 7 ( $\sum_{r=1}^3 C_3^r$ ) samples.

Given a caption text and several hallucinatory objects, we insert the objects into the appropriate locations in the caption, which can be effectively achieved with the help of GPT4. Automatically, the caption and objects are fed to the GPT4, with the prompt as *Add\_Prompt* (see Table 13).

**Removing existing Objects** Except from inserting hallucinatory objects, we also remove objects from the captions to create negative samples. We randomly select 1 or 2 segmented objects in the image which results in 6 negative samples ( $\sum_{r=1}^2 C_3^r$ ), and ask GPT4 to remove them from the caption with the *Remove\_Object\_Prompt*. To account for scenarios where the identified objects are not present in the title text, we ask GPT to alter elements like objects, colors, and properties in the original caption, the prompt we use is *Alter\_Object\_Prompt*. The prompt can be found in Table 13.

we construct a dataset of 500 samples for each of the COCO (Lin et al., 2014), Flickr30K (Young et al., 2014), and the out of domain subset of NoCaps Validation datasets (Agrawal et al., 2019a), with 27 negative samples for each image. Specifically, the out of domain subset of NoCaps comprises objects not seen in the COCO dataset, commonly used to measure a model’s ability to generalize to unseen classes.<sup>2</sup> The average length of the captions in the datasets is shown in Table 10.

<sup>2</sup>Our selection of Nocaps as the out-of-domain dataset is specific to our fine-tuning process in Section 4 and not the pre-training process of CLIP.

### 3.2 Evaluation and Analysis

We study several models to evaluate their performance on our benchmark. Each image is paired with a correct caption and 27 negative samples, and models are required to calculate the similarity between the image and the caption candidates and select the correct caption.

**Models** We evaluate a variety of models on our benchmark, including CLIP (Radford et al., 2021) ViT-B/32 and ViT-L/14; MetaCLIP (Xu et al., 2023) and DFN2B CLIP (Fang et al., 2023) are models pretrained on high-quality dataset after data curation; CLIPA (Li et al., 2023b) which achieves efficient training by using shorter image/text sequences, which reduces the computational load during the training period; EVA CLIP (Sun et al., 2023a) which employs innovative representation learning technology, optimization methods, and enhancement strategies to improve model performance; SigLIP (Zhai et al., 2023) which employs a contrastive learning loss function based on the Sigmoid function instead of the traditional softmax for pre-training on language and image data; CLIP ConvNext (Liu et al., 2022) is a variant of the CLIP model that uses ConvNext as the image encoder; CLIP NLLB-SigLip (Visheratin, 2023) is another variant that combines a text encoder from the NLLB model (Costa-jussà et al., 2022) and an image encoder from the SigLIP model; NegCLIP (Yüksekönül et al., 2023), an improved model based on CLIP ViT-B/32, which enhances the understanding of relationships between objects, attributes, and the sequence of words by swapping phrases; CECLIP (Zhang et al., 2023) which further develop enhanced negative samples and employ contrastive loss to enhance compositional reasoning; FLAVA (Singh et al., 2022) which is a single unified foundation model which can work across vision, language as well as vision-and-language multi-modal tasks; CoCa (Yu et al., 2022) is a pretrained model with contrastive and generative learning objectives; XVLM (Zeng et al., 2021) which aligns the visual concept and textual input in a multi-grained manner with 14M and 16M pretrained images; BLIP (Li et al., 2022) which effectively utilizes the noisy web data by bootstrapping the captions with 14M and 129M pretrained images; BLIP2 (Li et al., 2023a)<sup>3</sup> which bridges the gap between the visual and textual modalities with

<sup>3</sup>We use the image-text matching head for both BLIP and BLIP2.

Model	Params	OHD-Caps Benchmark			
		COCO	Flickr30K	NoCaps	Avg.
(a) comparisons with CLIP Models					
CLIP ViT-B/16	149M	16.6	17.2	8.6	14.1
CLIP ViT-B/32	151M	15.2	17.6	10.2	14.3
CLIP ViT-L/14	428M	22.4	22.6	12.0	19.0
MetaCLIP B/32	151M	25.6	25.2	16.0	22.3
MetaCLIP L/14	428M	36.8	26.4	19.4	27.5
CLIPA V2 L/16	428M	35.6	31.0	18.8	28.5
EVA-02 CLIP-B/16	149M	26.4	25.4	18.6	23.5
EVA-02 CLIP-L/14	428M	38.8	31.6	21.4	30.6
DFN2B CLIP B/16	149M	29.4	27.8	17.0	24.7
DFN2B CLIP L/14	428M	37.6	37.8	23.2	32.9
CLIP ConvNext-B	180M	34.0	28.0	20.4	27.5
CLIP ConvNext-L	352M	43.4	35.8	25.0	34.7
SigLIP B/16	203M	34.2	32.2	23.8	30.1
SigLIP L/16	652M	48.4	38.4	<b>30.8</b>	39.2
SigLIP SoViT-400m	877M	50.8	<b>41.4</b>	26.6	<b>39.6</b>
CLIP NLLB-SigLip-B	508M	25.2	20.0	22.6	22.6
CLIP NLLB-SigLip-L	1.1B	32.6	29.0	26.4	29.3
NegCLIP	151M	32.8	28.0	25.0	28.6
CECLIP	151M	<b>52.8</b>	40.8	23.4	39.0
(b) comparisons with other Image-Text Matching Models					
FLAVA	350M	28.0	28.4	16.6	24.3
CoCa	2.1B	26.0	24.4	20.0	23.5
XVLM 4M	216M	46.4	35.8	34.0	38.7
XVLM 16M	216M	41.8	19.4	21.8	27.7
BLIP 14M	583M	51.4	<b>48.0</b>	<b>42.0</b>	47.1
BLIP 129M	583M	40.8	38.0	31.2	36.7
BLIP2	3.4B	<b>62.6</b>	42.2	41.2	<b>48.7</b>

Table 1: Results of various models on our benchmark. NoCaps subset is used to evaluate zero-shot generalization.

a Q-former.

**Results** Table 1 shows the results of the models on our benchmark. From the results, we could find that,

- First of all, the vanilla CLIP models perform poorly across all three datasets, indicating their limited ability to recognize illusory objects in images. Multiple variants of CLIP, through improvements in data (e.g., MetaCLIP, DFN2B CLIP), model architecture (e.g., CLIP ConvNext, CLIP NLLB-SigLip), and training methods (e.g., CLIPA, EVA CLIP, SigLip), achieve a slight enhancement in the performance of the original CLIP. Among these variants, SigLIP demonstrates the most notable performance, exhibiting the best results on out-of-domain datasets and showcasing superior generalization capabilities.
- Secondly, NegCLIP attempts to enhance the model’s understanding of the text by parsing and substituting phrases, but it only achieves a marginal improvement compared to the original CLIP model. CECLIP exhibits relatively better

performance, which is mainly due to the constructed negative samples enhancing the model’s comprehension of the combined semantics of sentences. The NegCLIP and CECLIP models are trained on the COCO training set to distinguish between positive samples and enhanced negative samples. This might contribute to CECLIP’s good performance on the COCO dataset, owing in part to the model’s memory of the original correct text. However, their performance on the NoCaps dataset indicates that these models cannot effectively differentiate hallucinated objects.

- Furthermore, generative vision-language models typically achieve higher performance than vanilla CLIP models due to their more precise alignment of image and text representations. Furthermore, it is generally observed that the larger the model parameters, the better the performance. In particular, BLIP2, which has the highest number of parameters, performs best across all three datasets. In comparison, the XVLM 4M model has relatively fewer parameters but still demonstrates good performance. This indicates that XVLM’s strategy of multi-scale alignment indeed assists the model in more accurately capturing the fine-grained details within images.
- Finally, the overall trend among different models is consistent across the three datasets, with their performance typically being the lowest on the NoCaps dataset. Although fewer objects are recognized on the NoCaps dataset than Flickr30K, the performance is the lowest there due to the inclusion of categories that are out-of-domain. The BLIP 14M model demonstrates the best performance on both Flickr and NoCaps, which indicates its strong generalization capabilities.

**Analysis** The inability of models to recognize hallucinated objects primarily stems from the data used and the learning methods employed. The vanilla CLIP model is trained with a large number of image-caption pairs collected from the internet, using a contrastive loss function for optimization. Those captions are often brief and noisy, and the model is optimized to differentiate between correct and a multitude of incorrect image-text pairs. However, because the incorrect pairs are usually significantly different from the correct ones, the model can easily distinguish them. This means that the model does not need to learn the rich details in

the pictures to make accurate predictions. To address this issue, we need to make improvements to the original CLIP model in terms of data utilization and learning methodologies.

## 4 Methodology

We first revisit the training process of the vanilla CLIP model. Let  $I$  be the image and  $T$  be the text, the training objective of CLIP is to maximize the similarity between the image and text pairs and minimize the similarity between the image and text pairs that are not matched. The loss function is defined as:

$$\begin{aligned}\mathcal{L}_{i2t} &= -\log \frac{\exp(I \cdot T^+/\tau)}{\sum_{T^* \in \{T^+, T^-\}} \exp(I \cdot T^*/\tau)}, \\ \mathcal{L}_{t2i} &= -\log \frac{\exp(T \cdot I^+/\tau)}{\sum_{I^* \in \{I^+, I^-\}} \exp(T \cdot I^*/\tau)}, \\ \mathcal{L}_0 &= \frac{1}{2}(\mathcal{L}_{i2t} + \mathcal{L}_{t2i}),\end{aligned}\quad (1)$$

where  $T^+$  and  $I^+$  are the correct text and image, and  $T^-$  and  $I^-$  are the incorrect text and image, respectively.

With the addition of the negative samples  $T^{neg}$  created as in the previous section, we could modify the loss  $\mathcal{L}_{i2t}$  as:

$$\mathcal{L}_{i2t} = -\log \frac{\exp(I \cdot T^+/\tau)}{\sum_{T^* \in \{T^-, T^{neg}, T^+\}} \exp(I \cdot T^*/\tau)}.\quad (2)$$

To further enhance the model’s ability to distinguish between positive and negative samples, we additionally introduce a margin loss. This is to ensure that the distance between an image and its corresponding correct text is smaller than the distance to incorrect text by a specific threshold. This concept can be formulated as:

$$\mathcal{L}_1 = \max(0, \tau_1 - I \cdot T^+ + I \cdot T^*),\quad (3)$$

where  $\tau_1$  is the margin threshold, and  $T^* = \{T^-, T^{neg}\}$ .

Additionally, we generate enhanced negative samples by introducing perturbations to the original positive samples. Such negative samples are typically more challenging to distinguish than other negative samples within the batch. To encourage the model to recognize the partially correct information contained in the enhanced negative samples, resulting in a higher similarity to the positive samples compared to other negative samples within the

Model	OHD-Caps			
	COCO	Flickr30k	NoCaps	Avg.
Random	3.6	3.6	3.6	3.6
(a) comparisons with CLIP-Base baselines				
CLIP-B/32	15.2	17.6	10.2	14.3
NegCLIP	32.8	28.0	25.0	28.6
CECLIP	52.8	40.8	23.4	39.0
<b>Ours-B/32</b>	<b>80.4</b>	<b>85.0</b>	<b>82.0</b>	<b>82.5</b>
(b) comparisons with CLIP-Large baselines				
CLIP-L/14	26.0	27.0	16.8	23.3
<b>Ours-L/14</b>	<b>87.0</b>	<b>91.0</b>	<b>88.4</b>	<b>88.8</b>

Table 2: Results on OHD-Caps. CLIP-B/32, and CLIP-L/14 represent CLIP ViT-B/32 and CLIP ViT-L/14 336 px respectively.

batch, we introduce a margin loss between the in-batch negative samples and the enhanced negative samples:

$$\mathcal{L}_2 = \max(0, \tau_2 - I \cdot T^{neg} + I \cdot T^-), \quad (4)$$

where  $\tau_2$  is the margin threshold.

Next, we assign different weights to the aforementioned loss terms, allowing the model to learn adaptively. Consequently, the final loss function can be expressed as follows:

$$\mathcal{L} = \frac{1}{2}(\mathcal{L}_{t2i} + \mathcal{L}_{i2t}) + \lambda_1 \mathcal{L}_1 + \lambda_2 \mathcal{L}_2. \quad (5)$$

## 5 Experiments

**Training Datasets** We sample 8k images from the training set of COCO and 8k images from Flickr30k datasets, then generate negative samples for each image as in Section 3. Additionally, we randomly select  $\sim 1k$  samples from the COCO dataset’s validation set as our dev set for the selection of hyper-parameters. Detailed information about the dataset is provided in Table 10.

**Training Details** We utilize the CLIP ViT/32-B and CLIP ViT/14-L-336px implemented by Huggingface (Wolf et al., 2020) as the initial models and conduct fine-tuning for 10 epochs. The training process is carried out on a single A6000 GPU, with batch sizes of 56 and 14 set for the base and large models, respectively, and the learning rate is set at  $1e-6$ . The selection of hyper-parameters is determined by their performance on the validation set, where  $\lambda_1$  and  $\lambda_2$  are set as 0.1 and 0.1,  $\tau_1$  and  $\tau_2$  are set as 2.

	CIFAR-10 (2009)	CIFAR-100 (2009)	ImageNet (2009)	Eurosat (2019)	GTSRB (2012)	STL10 (2011)	avg. top-1 acc.
(a) comparisons with CLIP-Base baselines							
CLIP-B/32	<b>89.8</b>	64.2	<b>63.3</b>	46.3	<b>32.6</b>	<b>97.1</b>	65.6
NegCLIP	85.9	60.9	55.7	31.9	26.8	95.8	55.8
CECLIP	81.1	55.0	40.4	41.9	20.6	95.6	59.5
<b>Ours-B/32</b>	89.1	<b>66.0</b>	60.5	<b>51.7</b>	31.9	96.5	<b>66.0</b>
(b) comparisons with CLIP-Large baselines							
CLIP-L/14	95.0	74.4	<b>76.6</b>	61.4	<b>52.4</b>	99.4	<b>76.5</b>
<b>Ours-L/14</b>	95.0	<b>74.8</b>	72.8	<b>67.3</b>	43.6	99.4	75.5

Table 3: Zero-shot results on various datasets. The last column displays the average performance across 7 datasets.

**Evaluation** To verify the impact of our method on the model’s generalization capabilities, we conducted zero-shot experiments on the following datasets: CIFAR-10/100 (Krizhevsky et al., 2009), ImageNet-1K (Deng et al., 2009), DTD (Cimpoi et al., 2014), Eurosat (Helber et al., 2019), GTSRB (Stallkamp et al., 2012) and STL10 (Coates et al., 2011).

### 5.1 Main Results

We present the results for our self-constructed dataset in Table 2, and various zero-shot datasets in Table 3. From the results, we could find:

- Our model shows comparable zero-shot performance to vanilla CLIP Models (65.6 vs 66.0) and achieves significant improvements in hallucination recognition (14.3 vs 82.5). NegCLIP and CECLIP enhance the model’s capability of understanding composites by constructing negative samples and also achieve a moderate improvement on the OHD-Caps benchmark, with performance rising from 14.3% to 39.0%. However, the zero-shot performance of NegCLIP and CECLIP significantly decreases. This could be due to their reliance on rule-based methods to construct negative samples (such as swapping phrases), which may interfere with the model’s understanding of sentence semantics.
- Our model also demonstrates strong generalization capabilities in hallucination recognition. NegCLIP, CECLIP, and our model are all fine-tuned on the training set of the COCO dataset. Although they show varying degrees of performance improvement in COCO-related hallucination tests (NegCLIP at 32.8%, CECLIP at

Dataset	Criterion	Full Fine FT		LoRA FT	
		LLaVA	Ours	LLaVA	Ours
COCO	Accuracy ( $\uparrow$ )	<b>85.4</b>	81.2	85.7	<b>88.3</b>
	Precision ( $\uparrow$ )	81.8	<b>90.9</b>	81.8	<b>89.7</b>
	Recall ( $\uparrow$ )	<b>91.9</b>	85.1	<b>92.5</b>	86.9
	F1 Score ( $\uparrow$ )	86.4	<b>87.9</b>	86.7	<b>88.2</b>
	Yes ( $\rightarrow$ 50%)	56.5	<b>46.9</b>	56.8	<b>48.6</b>
Flickr30K	Accuracy ( $\uparrow$ )	73.7	<b>81.2</b>	74.4	<b>82.8</b>
	Precision ( $\uparrow$ )	67.5	<b>78.5</b>	67.9	<b>83.0</b>
	Recall ( $\uparrow$ )	<b>96.9</b>	88.0	<b>96.9</b>	85.7
	F1 Score ( $\uparrow$ )	79.2	<b>82.7</b>	79.5	<b>83.5</b>
	Yes ( $\rightarrow$ 50%)	73.1	<b>56.8</b>	72.5	<b>52.9</b>
NoCaps	Accuracy ( $\uparrow$ )	76.7	<b>81.3</b>	76.7	<b>82.6</b>
	Precision ( $\uparrow$ )	71.2	<b>80.6</b>	71.2	<b>81.8</b>
	Recall ( $\uparrow$ )	<b>92.7</b>	84.0	<b>92.3</b>	84.9
	F1 Score ( $\uparrow$ )	80.2	<b>82.0</b>	80.2	<b>83.2</b>
	Yes ( $\rightarrow$ 50%)	66.0	<b>52.7</b>	65.6	<b>52.3</b>

Table 4: Results on expanded POPE datasets. Yes denotes the proportion of answering ‘‘Yes’’ to the given question.

Model	Full FT				LoRA FT			
	$C_S \downarrow$	$C_I \downarrow$	Cover $\uparrow$	Length	$C_S \downarrow$	$C_I \downarrow$	Cover $\uparrow$	Length
LLaVA	56.4	14.9	79.1	106.4	58.2	16.4	<b>79.9</b>	106.5
<b>Ours</b>	<b>55.0</b>	<b>14.5</b>	<b>79.2</b>	107.5	<b>56.8</b>	<b>14.9</b>	79.2	108.5

Table 5: CHAIR hallucination evaluation results (max new tokens is 512) on COCO dev set. Smaller values correspond to less hallucinations.

52.8%), their performances are worse when facing unknown categories (NegCLIP at 25.0%, CE-CLIP at 23.4% for NoCaps images), indicating limited generalization capabilities of the models. In contrast, our model performs consistently across three different datasets, at approximately 82%. This result verifies that our model can effectively distinguish hallucinated objects in different datasets and possesses the capability to generalize across datasets.

## 5.2 Evaluation for LLaVA

To verify the effectiveness of the enhanced CLIP model compared to the original CLIP in assisting large vision-language models to mitigate the issue of object hallucination, we replace the CLIP ViT-L/14-336px baseline model in LLaVA-1.5 with our fine-tuned version. We train LLaVA (Liu et al., 2023c) from scratch using the hyper-parameters specified in the original paper. Comparison results with other methods, such as constructing SFT data (Wang et al., 2023a) or introducing DPO processes (Zhou et al., 2024; Zhao et al., 2023) for

Dataset	Criterion	Full FT		LoRA FT	
		LLaVA	Ours	LLaVA	Ours
Generative	$C_S \downarrow$	7.2	<b>6.5</b>	7.2	<b>6.1</b>
	$C_I \downarrow$	35.4	<b>31.7</b>	33.4	<b>30.1</b>
	Cover ( $\uparrow$ )	<b>52.2</b>	50.9	<b>51.7</b>	50.7
Discriminative	Accuracy ( $\uparrow$ )	74.3	<b>80.2</b>	74.2	<b>80.8</b>
	Precision ( $\uparrow$ )	<b>93.9</b>	85.5	<b>93.5</b>	86.4
	Recall ( $\uparrow$ )	65.6	<b>84.4</b>	65.7	<b>84.3</b>
	F1 ( $\uparrow$ )	77.2	<b>84.9</b>	77.2	<b>85.3</b>

Table 6: Results on AMBER dataset which includes the assessment of hallucinations in both discriminative and generative responses.

Model	Existence	Attribute	State	Number	Action	Relation
(a) Full FT						
LLaVA	83.5	72.4	67.0	78.7	85.2	57.4
<b>Ours</b>	<b>94.2</b>	<b>79.1</b>	<b>77.1</b>	<b>79.5</b>	<b>88.6</b>	<b>64.3</b>
(b) LoRA FT						
LLaVA	83.0	73.2	71.7	73.2	81.8	56.5
<b>Ours</b>	<b>94.3</b>	<b>79.4</b>	<b>77.8</b>	<b>80.4</b>	<b>86.7</b>	<b>63.4</b>

Table 7: Detailed performance on AMBER discriminative subset which includes evaluation results of other types of hallucinations, such as attribute, number, and relation.

further alignment can be found in Appendix B.

**Hallucination Detection** To evaluate the occurrence of hallucination phenomena in discriminative and generative responses within models, we select the following evaluation methods for analysis: an extended version of the POPE dataset (Li et al., 2023c) for discriminative response evaluation, and CHAIR evaluation (Rohrbach et al., 2018) for generative response; the AMBER dataset (Wang et al., 2023b) contains both types of evaluations. The format of the question contained in POPE is: ‘Is there a X in the image?’, where X refers to the name of the object. The questions in the dataset are designed such that the objects are present and absent in equal measure, therefore the ideal ‘yes’ response rate should be around 50%. We extend the POPE dataset and incorporate the Flickr30k and NoCaps domains to test the model’s generalization capabilities. The CHAIR metric evaluates object hallucinations in image descriptions by measuring the ratio of referenced objects not found in the ground-truth label set, with  $CHAIR_S$  for sentence level:

$$C_S = \frac{|\{ \text{hallucinated objects} \}|}{|\{ \text{all mentioned objects} \}|},$$

Model	MME	VQA-v2	VisWiz	SciQA-IMG	TextVQA
(a) Full FT					
LLaVA	1459.4	79.1	48.9	<b>69.4</b>	<b>58.5</b>
<b>Ours</b>	<b>1487.2</b>	<b>79.2</b>	<b>50.0</b>	69.3	58.2
(b) LoRA FT					
LLaVA	1445.4	79.1	46.8	<b>69.8</b>	<b>58.5</b>
<b>Ours</b>	<b>1455.4</b>	<b>79.2</b>	<b>47.2</b>	68	58.4

Table 8: Results on various benchmarks.

CHAIR<sub>I</sub> for image-level analysis:

$$C_I = \frac{|\{ \text{captions w/ hallucinated objects} \}|}{|\{ \text{all captions} \}|},$$

and Cover measures the object coverage of responses:

$$\text{Cover} = \frac{|\{ \text{captions w/ hallucinated objects} \}|}{|\{ \text{ground truth objects} \}|}.$$

Table 4, 5, 6 show the results of the expanded POPE dataset, CHAIR evaluation, and AMBER dataset, respectively. From the results, we could find:

- For discriminative responses, our model achieves significant improvements on various datasets. On the POPE dataset, compared to the original, it attains a better balance between accuracy and recall which results in a higher F1 score and also approaches a more ideal balance in the proportion of "Yes" responses. The same phenomenon of performance improvement is also observed in the AMBER dataset.
- For generative responses, our model demonstrates a lower proportion of hallucinated content on the COCO validation set and the AMBER dataset, while maintaining a relatively stable coverage and response length.

**General Performance** We evaluate the model’s general performance on different datasets, which include: MME-Perception (Fu et al., 2023) evaluates the model’s visual perception with yes/no questions. VQA-v2 (Goyal et al., 2017) evaluate model’s visual perception capabilities on open-ended short answers; VizWiz (Gurari et al., 2018) and ScienceQA (Lu et al., 2022) with multiple choice to evaluate the model’s zero-shot generalization on visual questions; TextVQA (Singh et al., 2019) contains text-rich visual question answering.

Model	$\mathcal{L}_0$	$\mathcal{L}_1$	$\mathcal{L}_2$	OHD-Caps	CIFAR10	CIFAR100	Avg.
CLIP				14.3	89.8	64.2	39.4
Ours	✓			80.1	88.6	<b>66.4</b>	79.1
	✓	✓		80.5	<b>89.3</b>	66.0	79.4
	✓		✓	81.6	89.0	66.3	80.0
	✓	✓	✓	<b>82.5</b>	89.1	66.0	<b>80.5</b>

Table 9: Ablation of losses on CLIP ViT-B/32.

Results are shown in Table 8. We can observe that with full fine-tuning, there is a slight improvement in the model’s average performance. Specifically, the average performance of the model across five datasets increased from 343.1 to 348.5, with the most notable improvement on the MME dataset. Conversely, when employing LoRA fine-tuning, the average performance of the model remained unchanged (340.0 vs 341.7).

### 5.3 Ablation Study

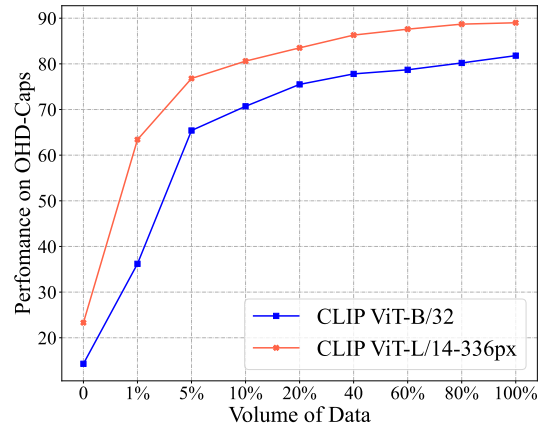


Figure 2: The performance of the model on the OHD-Caps dataset with different training data volumes provided. We report the average results of three random seeds.

In this subsection, we present ablation studies to examine the impact of our model’s different components. We conduct these experiments on the CLIP ViT-B/32 model.

**Losses** As demonstrated in Table 9, the inclusion of the  $\mathcal{L}_0$  loss alone significantly improves OHD-Caps performance over the baseline. Subsequently, iterative incorporation of  $\mathcal{L}_1$  and  $\mathcal{L}_2$  provide incremental benefits, with the full combination yielding the highest average performance. Compared to  $\mathcal{L}_1$  loss,  $\mathcal{L}_2$  loss has a more significant effect on improving model performance. This suggests that by increasing the distance between constructed negative samples and other negative samples in the



batch, the model can achieve a more refined understanding.

**Data Volume** Figure 2 shows the performance of the OHD-Caps dataset with varying amounts of training data. As can be seen from the figure, even with a very small amount of data, the model’s performance can be significantly improved. For example, by training with just 1% of the data (that is, 160 images), the performance of the CLIP-L/14 model can increase from 20% to 60%. However, as more data is added, the performance improvement gradually slows and stabilizes.

## 6 Conclusion

Our study investigates the reasons behind object hallucination in LVLMS. We construct a benchmark specifically for the evaluation of hallucinations and find that the visual perception module commonly used in current LVLMS, i.e., the CLIP model, cannot effectively discriminate hallucinated text. By designing negative samples and optimizing the contrastive loss function, we achieve a significant improvement in model performance on the hallucination detection dataset. Moreover, replacing the original CLIP model with our improved model can effectively alleviate the issue of object hallucination in the LLaVA model.

## Limitations

Although we conducted a series of explorations, our research still has its limitations. Firstly, our focus is solely on the issue of object hallucination within LVLMS, and we do not extend our research to other types of hallucinations. Secondly, the benchmark we propose comprises over 20 negative samples. Due to budgetary constraints, the size of this dataset is much smaller compared to the datasets used for evaluating compositional understanding, e.g. ARO dataset (Yüksekgönül et al., 2023). Thirdly, we only evaluate the visual encoders of most LVLMS, i.e. the CLIP models, but we do not conduct research on encoders used by some other models, for instance, the variant of ResNet called NFNet-F6 (Brock et al., 2021) used by Flamingo (Alayrac et al., 2022).

## Ethics Statement

Object hallucination severely limits the practical application of LVLMS. For example, in medical image diagnosis, it can lead to false descriptions

of tumor objects that are not present in the image. While our work has mitigated hallucinations in the visual encoder of LVLMS, hallucinations may still exist in the multi-head attention layers and feed-forward layers. Real-world applications based on LVLMS must systematically control hallucinations to avoid negative impacts on users.

## Acknowledgement

The authors wish to thank all reviewers for their helpful comments and suggestions. The corresponding authors are Yuanbin Wu and Aimin Zhou. This research was (partially) supported by NSFC(62076097), National Key R&D Program of China (2021YFC3340700), the Open Research Fund of Key Laboratory of Advanced Theory and Application in Statistics and Data Science (East China Normal University), Ministry of Education.

## References

- Harsh Agrawal, Peter Anderson, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019a. [nocaps: novel object captioning at scale](#). In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 8947–8956. IEEE.
- Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, and Peter Anderson. 2019b. [Nocaps: Novel object captioning at scale](#). In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8948–8957.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022. [Flamingo: a visual language model for few-shot learning](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. [Vqa: Visual question answering](#). In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.
- Ali Furkan Biten, Lluís Gómez, and Dimosthenis Karatzas. 2022. [Let there be a clock on the beach](#):

- [Reducing object hallucination in image captioning](#). In *2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2473–2482.
- Andy Brock, Soham De, Samuel L. Smith, and Karen Simonyan. 2021. [High-performance large-scale image recognition without normalization](#). In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 1059–1071. PMLR.
- 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](#).
- M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. 2014. Describing textures in the wild. In *CVPR*.
- Adam Coates, Andrew Ng, and Honglak Lee. 2011. An analysis of single-layer networks in unsupervised feature learning. In *AISTAT*.
- Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Y. Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loïc Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. [No language left behind: Scaling human-centered machine translation](#). *CoRR*, abs/2207.04672.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*.
- Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Sarkar, Rohan Taori, Yusuke Noda, Demetri Terzopoulos, Yejin Choi, et al. 2024. Agent ai: Surveying the horizons of multimodal interaction. *arXiv preprint arXiv:2401.03568*.
- Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander Toshev, and Vaishaal Shankar. 2023. [Data filtering networks](#). *CoRR*, abs/2309.17425.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, and Rongrong Ji. 2023. [MME: A comprehensive evaluation benchmark for multimodal large language models](#). *CoRR*, abs/2306.13394.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. [Making the V in VQA matter: Elevating the role of image understanding in visual question answering](#). In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 6325–6334. IEEE Computer Society.
- Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P. Bigham. 2018. [Vizwiz grand challenge: Answering visual questions from blind people](#). In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 3608–3617. Computer Vision Foundation / IEEE Computer Society.
- Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. 2019. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*
- Cheng-Yu Hsieh, Jieyu Zhang, Zixian Ma, Aniruddha Kembhavi, and Ranjay Krishna. 2023. [Sugarcrepe: Fixing hackable benchmarks for vision-language compositionality](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- Hongyu Hu, Jiyuan Zhang, Minyi Zhao, and Zhenbang Sun. 2023. [CIEM: contrastive instruction evaluation method for better instruction tuning](#). *CoRR*, abs/2309.02301.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. [Survey of hallucination in natural language generation](#). *ACM Comput. Surv.*, 55(12):248:1–248:38.
- Amita Kamath, Jack Hessel, and Kai-Wei Chang. 2023. What’s “up” with vision-language models? investigating their struggle with spatial reasoning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9161–9175.
- Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images.
- 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](#). *CoRR*, abs/2311.16922.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. 2023a. [BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models](#). In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 19730–19742. PMLR.

- Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. 2022. [BLIP: bootstrapping language-image pre-training for unified vision-language understanding and generation](#). In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 12888–12900. PMLR.
- Xianhang Li, Zeyu Wang, and Cihang Xie. 2023b. [Clipa-v2: Scaling CLIP training with 81.1% zero-shot imagenet accuracy within a \\$10, 000 budget; an extra \\$4, 000 unlocks 81.8% accuracy](#). *CoRR*, abs/2306.15658.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023c. [Evaluating object hallucination in large vision-language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 292–305. Association for Computational Linguistics.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. [Microsoft COCO: common objects in context](#). In *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V*, volume 8693 of *Lecture Notes in Computer Science*, pages 740–755. Springer.
- 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](#). *CoRR*, abs/2310.14566.
- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2023b. [Aligning large multi-modal model with robust instruction tuning](#). *CoRR*, abs/2306.14565.
- Hanchao Liu, Wenyan 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](#). *CoRR*, abs/2402.00253.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023c. [Visual instruction tuning](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. 2022. [A convnet for the 2020s](#). In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 11966–11976. IEEE.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. [Learn to explain: Multimodal reasoning via thought chains for science question answering](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. [Learning transferable visual models from natural language supervision](#). In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.
- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2018. [Object hallucination in image captioning](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 4035–4045. Association for Computational Linguistics.
- Ravi Shekhar, Sandro Pezzelle, Yauhen Klimovich, Aurélie Herbelot, Moin Nabi, Enver Sangineto, and Raffaella Bernardi. 2017. [FOIL it! find one mismatch between image and language caption](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 255–265. Association for Computational Linguistics.
- Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. 2022. [FLAVA: A foundational language and vision alignment model](#). In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 15617–15629. IEEE.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. [Towards VQA models that can read](#). In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 8317–8326. Computer Vision Foundation / IEEE.
- Johannes Stallkamp, Marc Schlipf, Jan Salmen, and Christian Igel. 2012. [Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition](#). *Neural networks*.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F. Christiano. 2020. [Learning to summarize with human feedback](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing*

- Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.*
- Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. 2023a. [EVA-CLIP: improved training techniques for CLIP at scale](#). *CoRR*, abs/2303.15389.
- Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, Kurt Keutzer, and Trevor Darrell. 2023b. [Aligning large multimodal models with factually augmented RLHF](#). *CoRR*, abs/2309.14525.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#). *CoRR*, abs/2302.13971.
- Alexander A. Visheratin. 2023. [NLLB-CLIP - train performant multilingual image retrieval model on a budget](#). *CoRR*, abs/2309.01859.
- Junke Wang, Lingchen Meng, Zejia Weng, Bo He, Zuxuan Wu, and Yu-Gang Jiang. 2023a. [To see is to believe: Prompting GPT-4V for better visual instruction tuning](#). *CoRR*, abs/2311.07574.
- Junyang Wang, Yuhang Wang, Guohai Xu, Jing Zhang, Yukai Gu, Haitao Jia, Ming Yan, Ji Zhang, and Jitao Sang. 2023b. [An llm-free multi-dimensional benchmark for mllms hallucination evaluation](#). *CoRR*, abs/2311.07397.
- Junyang Wang, Yiyang Zhou, Guohai Xu, Pengcheng Shi, Chenlin Zhao, Haiyang Xu, Qinghao Ye, Ming Yan, Ji Zhang, Jihua Zhu, Jitao Sang, and Haoyu Tang. 2023c. [Evaluation and analysis of hallucination in large vision-language models](#). *CoRR*, abs/2308.15126.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. 2023. [The rise and potential of large language model based agents: A survey](#). *arXiv preprint arXiv:2309.07864*.
- Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. 2023. [Demystifying CLIP data](#). *CoRR*, abs/2309.16671.
- Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. [From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions](#). *Trans. Assoc. Comput. Linguistics*, 2:67–78.
- Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. 2022. [Coca: Contrastive captioners are image-text foundation models](#). *Trans. Mach. Learn. Res.*, 2022.
- Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. 2016. [Modeling context in referring expressions](#). In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, pages 69–85. Springer.
- Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun, and Tat-Seng Chua. 2023. [RLHF-V: towards trustworthy mllms via behavior alignment from fine-grained correctional human feedback](#). *CoRR*, abs/2312.00849.
- Mert Yükekönül, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. 2023. [When and why vision-language models behave like bags-of-words, and what to do about it?](#) In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Yan Zeng, Xinsong Zhang, and Hang Li. 2021. [Multi-grained vision language pre-training: Aligning texts with visual concepts](#). *CoRR*, abs/2111.08276.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. [Sigmoid loss for language image pre-training](#). In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pages 11941–11952. IEEE.
- Le Zhang, Rabiul Awal, and Aishwarya Agrawal. 2023. [Contrasting intra-modal and ranking cross-modal hard negatives to enhance visio-linguistic fine-grained understanding](#). *CoRR*, abs/2306.08832.
- Yi-Fan Zhang, Weichen Yu, Qingsong Wen, Xue Wang, Zhang Zhang, Liang Wang, Rong Jin, and Tieniu Tan. 2024. [Debiasing multimodal large language models](#).
- Zhiyuan Zhao, Bin Wang, Linke Ouyang, Xiaoyi Dong, Jiaqi Wang, and Conghui He. 2023. [Beyond hallucinations: Enhancing llms through hallucination-aware direct preference optimization](#). *CoRR*, abs/2311.16839.
- Yiyang Zhou, Chenhang Cui, Rafael Rafailov, Chelsea Finn, and Huaxiu Yao. 2024. [Aligning modalities in vision large language models via preference fine-tuning](#). *CoRR*, abs/2402.11411.

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](#). *CoRR*, abs/2310.00754.

Xueyan Zou, Jianwei Yang, Hao Zhang, Feng Li, Linjie Li, Jianfeng Wang, Lijuan Wang, Jianfeng Gao, and Yong Jae Lee. 2023. [Segment everything everywhere all at once](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

## A Statistics on the Datasets

Dataset	Size	#Negative Samples	#Avg Length
<i>Train</i>			
COCO	8000	27	16.0
Flickr30K	8000	27	18.4
<i>Dev</i>			
COCO	990	27	15.6
<i>Test</i>			
COCO	500	27	16.3
Flickr30K	500	27	21.1
Nocaps	500	27	19.1

Table 10: Statistics of the datasets used in our benchmark.

The statistical information of the dataset is presented in the Table 10, which is divided into three parts: training, testing, and validation. The average length displayed in the table refers to the average length of the negative examples in the dataset.

## B Comparison with Other Methods

To demonstrate that the proposed method has fewer object hallucinations and better general performance than other popular methods, we additionally compared the following approaches: LVIS (Wang et al., 2023a) built a 220k visual instruction dataset. By utilizing the excellent visual analysis ability of GPT-4V and generating data through carefully designed prompts. Expanding the original LLaVA training data, datasets of different sizes, 619k and 880k, were obtained; POVID (Zhou et al., 2024) and DPO (Zhao et al., 2023) build hallucination texts using GPT4V and GPT4 respectively, and compose pairs with high-quality non-illusionary replies for DPO optimization. We report the model results based on the checkpoints provided by the paper.

The results are shown in Table 11 and Table 12. From the results, our method outperforms the instruction finetune-based and dpo-based methods in

Model	COCO		Flickr30K		Nocaps	
	F1	Yes %	F1	Yes %	F1	Yes %
(a) Full FT						
LLaVA	86.4	56.5	79.2	73.1	80.2	66.0
LVIS-619k	77.4	32.6	70.2	33.6	67.3	31.2
LVIS-880k	85.6	41.7	79.7	<b>45.6</b>	80.6	43.7
<b>Ours</b>	<b>87.9</b>	<b>46.9</b>	<b>82.7</b>	56.8	<b>82.0</b>	<b>52.7</b>
(a) LoRA FT						
LLaVA	86.7	56.8	79.5	72.5	80.2	65.6
POVID	86.8	44.9	81.9	<b>51.8</b>	81.4	49.6
HADPO	84.6	43.0	75.1	43.5	78.4	43.7
<b>Ours</b>	<b>88.2</b>	<b>48.6</b>	<b>83.5</b>	52.9	<b>83.2</b>	<b>52.3</b>

Table 11: Comparison results on expanded POPE datasets. Yes% denotes the proportion of answering “Yes” to the given question.

Model	MME	VQAv2	VisWiz	SciQA-IMG	TextVQA
(a) Full FT					
LLaVA	1459.4	79.1	48.9	<b>69.4</b>	58.5
LVIS-619k	1473.6	79.2	50.0	68.1	57.7
LVIS-880k	1517.7	<b>79.6</b>	<b>51.7</b>	68.9	<b>58.7</b>
<b>Ours</b>	<b>1487.2</b>	79.2	50.0	69.3	58.2
(b) LoRA FT					
LLaVA	1445.4	79.1	46.8	69.8	<b>58.5</b>
POVID	1418.5	78.8	42.3	67.5	58.0
HADPO	1430.4	76.4	43.4	<b>70.3</b>	56.6
<b>Ours</b>	<b>1455.4</b>	<b>79.2</b>	<b>47.2</b>	68	58.4

Table 12: Comparison Results on various benchmarks.

terms of performance on POPE (our method improved the average F1 score by 2.6, while LVIS, HADPO, and POVID showed no significant improvement), demonstrating lower hallucination rates. Additionally, our method shows comparable performance to other methods in terms of general performance.

## C More Examples

We present more examples in Figure 3. It can be observed that our method can seamlessly integrate objects that are not present in the original image into the text. The names of the added objects are highlighted in red. Removing objects that are present in the picture can be accomplished with minimal adjustments. As for the removal of objects not depicted in the image, such as the “food” mentioned in the third figure, the negative samples typically involve modifications to the objects, attributes, and other content in the positive samples.

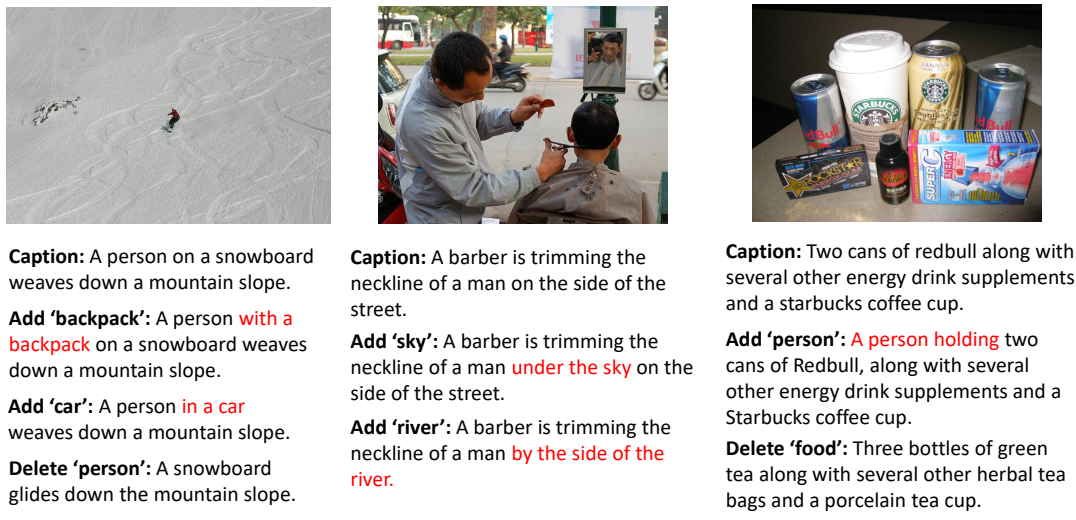


Figure 3: Examples from our benchmark OHD-Caps. The three images in the figure are from the COCO, Flickr, and Nocaps datasets, respectively.

<p><b>Prompt Template</b></p> <hr/> <p><b>Add_Prompt:</b> <i>Given a sentence {caption}, generate a new sentence and includes each object from the list {objects}. Make the changes to the original sentence as minimal as possible. Ensure that the new sentence is coherent, natural, semantically smooth and free of grammatical errors.</i></p> <hr/> <p><b>Remove_Object_Prompt:</b> <i>Given a sentence {caption}, generate a new sentence and remove each object from list {objects} to make the semantics of the sentence different. Ensure that the new sentence is coherent, natural, semantically smooth and free of grammatical errors.</i></p> <hr/> <p><b>Alter_Object_Prompt:</b> <i>Given a sentence {caption}, choose to modify the objects, colors, attributes, etc., within the sentence to make the semantics of the sentence different. Make the changes to the original sentence as minimal as possible. Ensure that the new sentence is coherent, natural, semantically smooth and free of grammatical errors.</i></p>
---

Table 13: Prompt Templates for Querying GPT-4. We replace the object that is to be added or deleted with object in the prompt, and replace caption with the original caption text. The revised text should then be submitted to GPT-4 to generate the corresponding output.

**D Prompt Template**

Table 13 presents the prompt templates for generating negative samples that we used in Section 3.