# Mitigating Object Hallucination in Large Vision-Language Models via Classifier-Free Guidance

Linxi Zhao[*†]  Yihe Deng[*‡]  Weitong Zhang[§]  Quanquan Gu[¶]

**Abstract**

The advancement of Large Vision-Language Models (LVLMs) has increasingly highlighted the critical issue of their tendency to hallucinate non-existing objects in the images. To address this issue, previous works focused on using specially curated datasets or powerful LLMs (e.g., GPT-3.5) to rectify the outputs of LVLMs. However, these approaches require either expensive training/fine-tuning or API access to advanced LLMs to correct the model's output post-generation. In this paper, we tackle this challenge by introducing a framework called **M**itigating hallucin**A**tion via classifie**R**-Free gu**I**da**N**c**E** (`MARINE`), which is both *training-free* and *API-free*, and can effectively and efficiently reduce object hallucinations during the generation process. Specifically, `MARINE` enriches the visual context of LVLMs by integrating existing open-source vision models, and employs classifier-free guidance to incorporate the additional object grounding features to improve the precision of LVLMs' generations. Through comprehensive evaluations across 6 popular LVLMs with diverse evaluation metrics, we demonstrate the effectiveness of `MARINE`, which even outperforms existing fine-tuning-based methods. Remarkably, it not only reduces hallucinations but also improves the detailedness of LVLMs' generations, as assessed by GPT-4V.

## 1 Introduction

The advent of Large Language Models (LLMs) has motivated advancements in extending their remarkable capabilities to multimodal data. Grounded in the development of pre-trained vision-language models (Radford et al., 2021; Jia et al., 2021; Alayrac et al., 2022) that align visual and textual embedding spaces, Large Vision Language Models (LVLMs) have gained substantial attention in both architectural development (Liu et al., 2023d; Zhu et al., 2023; Ye et al., 2023; Dai et al., 2023a; Gao et al., 2023) and benchmarking datasets (Xu et al., 2023; Lu et al., 2024). However, similar to the hallucination issues in textual LLMs (Ji et al., 2023), where irrelevant content is generated with input prompts, LVLMs face a specific challenge known as object hallucination: generating descriptions of non-existing objects for a given image (Li et al., 2023b; Wang et al., 2023b; Zhou et al., 2023; Fu et al., 2023; Lovenia et al., 2023). Such a problem is particularly concerning as it

---

[*]Equal contribution

[†]Xinya College, Tsinghua University, Haidian District, Beijing, 100084, P. R. China; e-mail: `zhaolx19@mails.tsinghua.edu.cn`

[‡]Department of Computer Science, University of California, Los Angeles, CA 90095, USA; e-mail: `yihedeng@cs.ucla.edu`

[§]Department of Computer Science, University of California, Los Angeles, CA 90095, USA; e-mail: `weightzero@ucla.edu`

[¶]Department of Computer Science, University of California, Los Angeles, CA 90095, USA; e-mail: `qgu@cs.ucla.edu`

compromises the model's accuracy and reliability, especially considering the growing application of LVLMs to safety-critical downstream tasks such as medical imaging (Chambon et al., 2022; Bazi et al., 2023).
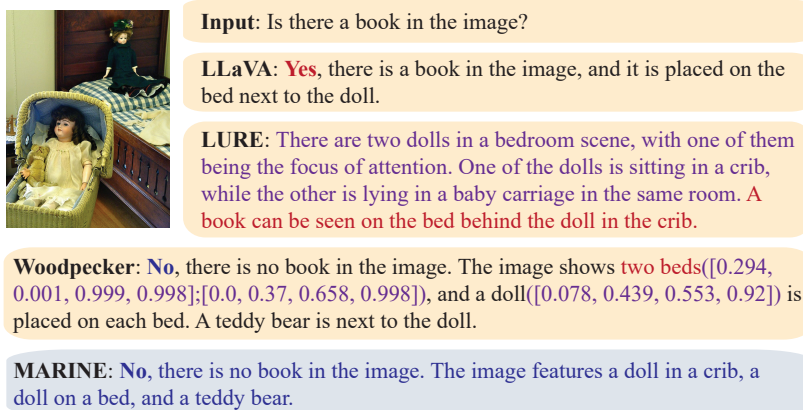


Figure 1: Example responses to an image-question pair. The LURE-corrected output deviates from the original question, offering irrelevant descriptions without directly addressing the query. Woodpecker similarly overwrites the original output, introducing bounding boxes which may not be of user's inquiry. It also hallucinates the existence of two beds while there is only one bed in the figure. In contrast, MARINE maintains the original answer's style and adheres to the user's instruction while eliminating hallucination.

In response to the pressing issue of object hallucinations in LVLMs, early attempts (Liu et al., 2023a,b; Gunjal et al., 2023; Wang et al., 2023a) have focused on correcting the bias inherent from the pre-training data by curating high-quality datasets for fine-tuning, specifically designed to mitigate object hallucinations. However, the creation of such extensive, high-quality datasets and the subsequent fine-tuning of LVLMs entail significant costs in terms of human annotation and computational resources. Consequently, recent works propose more cost-efficient strategies, employing post-generation correction methods with either a minimally fine-tuned corrector model (Zhou et al., 2023; Zhai et al., 2023) or advanced GPT APIs (Yin et al., 2023). While post-generation approaches effectively correct errors in the generated content, it is important to note that they also overwrite the original outputs of the models. This phenomenon is illustrated in Figure 1. In particular, the specialized input for post-generation (e.g., 'correct the hallucinated objects in this description.') can negatively impact the LVLM's inherent diversity in responding to different types of questions and adherence to the instruction. The corrector models, such as GPT-3.5, would also introduce inherent hallucinations on their own. Fine-tuning methods like LURE (Zhou et al., 2023) may further result in over-fitting of the annotations in the fine-tuning dataset.

To strike a balance between reducing object hallucinations, computational efficiency, limited/no access to advanced LLMs, and preserving LLM originality, we introduce **M**itigating hallucin**A**tion via classifie**R**-Free gu**I**da**N**c**E** (MARINE), a training-free, API-free[1] framework that performs corrections during the generation process. As shown in Figure 2, our framework incorporates a pre-trained object grounding vision encoder to enrich the visual context of LVLMs and controls the text generation

---

[1]The term "API-free" in denotes the elimination of any need for API calls to OpenAI. We note that Woodpecker requires 3-5k input tokens for an API call to each short captioning task.

via classifier-free guidance (CFG) (Ho and Salimans, 2021) specifically designed for the *multi-modal* setting. Extracting the visual features from the object grounding encoder and projecting it as a soft prompt to the LVLMs, we utilize CFG to generate the guided output, which places more importance on the enriched visual features from the object grounding encoder. It is important to highlight that our framework is compatible with any vision model and projection function. In our paper, we present results based on the DEtection TRansformer (DETR) (Carion et al., 2020), denoted by MARINE-DETR, as well as the ideal results based on ground truth object oracle, denoted by MARINE-Truth.

Empirical evaluations are conducted on six widely-recognized LVLMs using the MSCOCO dataset (Lin et al., 2014). Our experimental results demonstrate that, in comparison with state-of-the-art algorithms, MARINE exhibits further reduced hallucination, as measured by existing hallucination metrics such as CHAIR (Rohrbach et al., 2018) and POPE (Li et al., 2023b), as well as additional metrics considered in this study including the recall and GPT-4V evaluation on the detailedness of the responses. These promising results confirm that MARINE can effectively mitigate object hallucinations without requiring additional training resources or access to advanced LLMs. Moreover, our ablation studies elucidate the impact of different levels of guidance strength on performance. We also provide specific examples to illustrate how the guidance influences the output logits of the LVLMs.

**Notation.** We use lower case letters, lower case bold face letters, and upper case bold face letters to denote scalars, vectors, and matrices respectively. We use the symbol $p$ to represent the conditional probability of LLM's response. And we denote the sequence of tokens generated before the $t$-th token as $\mathbf{y}_{<t} = [y_1, \ldots, y_{t-1}]$ for $t > 1$.
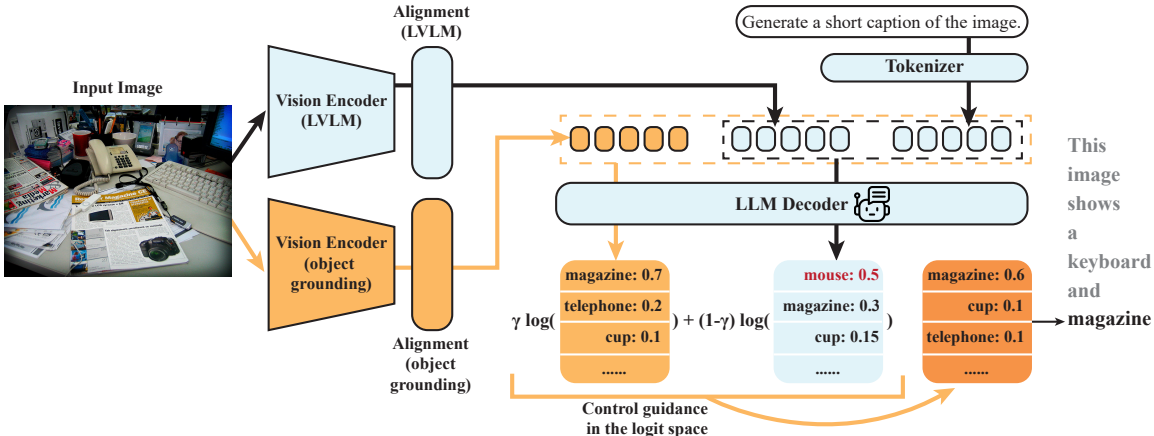


Figure 2: Illustration of MARINE framework, which adds an object grounding encoder with direct alignment to enrich the visual context of the original LVLM. The output logits are controlled to place more importance on this object grounding soft prompt with the guidance strength $\gamma$.

## 2 Related Work

### 2.1 Hallucination in Large Vision-Language Models

Since the introduction of recent Large Vision-Language Models (LVLMs) (Liu et al., 2023d; Zhu et al., 2023; Ye et al., 2023; Dai et al., 2023a; Gao et al., 2023), the hallucination phenomenon in these models has gathered significant attention in the research community. This issue was first highlighted by Li et al. (2023b) with subsequent studies (Wang et al., 2023b; Zhou et al., 2023; Fu

et al., 2023; Lovenia et al., 2023) that, LVLMs exhibit similar hallucination problems as the textual LLMs. Notably, different from textual LLMs, LVLMs are prone to a unique type of hallucination called 'object hallucination' (Rohrbach et al., 2018), where the model falsely perceives the presence of non-existent objects in images.

In response to object hallucination problems, efforts have been made to mitigate object hallucination in smaller image captioning models (Biten et al., 2022; Dai et al., 2023b). Regarding the recent development of LVLMs, several works (Liu et al., 2023b; Gunjal et al., 2023) proposed vision-language fine-tuning datasets aimed for improved robustness. Wang et al. (2023a) leveraged the vision-language model to generate more diverse instruction-tuning data and iteratively correct the inaccuracies in data. Zhai et al. (2023) introduced a GPT-4 assisted evaluation method and also a fine-tuning strategy using the MSCOCO dataset. This fine-tuning approach incorporates a binary switching parameter $\epsilon \in \{\pm 1\}$ within the linear projection layer, trained according to the discrepancy between the outputs of an object detector and the ground truth values. Most related to our setting, Yin et al. (2023) proposed Woodepecker, a five-stage training-free method eventually leveraging GPT-3.5 API for hallucination correction. Concurrently, Leng et al. (2023) proposed Visual Contrastive Decoding (VCD), which involves distorting image inputs with noise and imposing penalties on the logit outputs of these corrupted images. Instead of corrupting the image and adding penalty, our method introduces additional visual features to guide the generation.

## 2.2 Controllable Generation

Controllable text generation (Prabhumoye et al., 2020; Hu and Li, 2021; Zhang et al., 2023a) has emerged as a vital research domain, focusing on the generation of natural sentences with controllable attributes such as persona (Prabhumoye et al., 2020; Hu and Li, 2021; Zhang et al., 2023a), politeness (Niu and Bansal, 2018; Madaan et al., 2020), and story ending (Peng et al., 2018). Among the various approaches, fine-tuning has been recognized as the most straightforward approach, achieved either through the tuning of model parameters (Li and Liang, 2021; Ouyang et al., 2022; Carlsson et al., 2022) or the integration of tunable adaptor modules (Lin et al., 2021; Ribeiro et al., 2021). While fine-tuning has been effective in a wide range of applications, it is also expensive in computational cost as the size of LLMs is growing tremendously. Recently, there has been a development on controllable generation with diffusion models (Li et al., 2022; Lin et al., 2023), extending to controllable text-to-image generation (Yang et al., 2023). Particularly, the use of classifier guidance (Dhariwal and Nichol, 2021) and classifier-free guidance (Ho and Salimans, 2021) has become prominent in refining the quality of generated outputs. While classifier guidance employs an auxiliary classifier model to evaluate and improve the generation Kawar et al. (2022); Kim et al. (2022); Shi et al. (2023), classifier-free guidance integrates control directly into the generative model, offering an efficient approach for real-time applications with computational constraints (Saharia et al., 2022; Lin et al., 2024). Most recently, Sanchez et al. (2023) applied classifier-free guidance to language models in the *single-modal* setting to improve their performance at inference time. Inspired by these recent developments, we introduce a novel methodology in the *multi-modal* setting aimed at reducing hallucinations with classifier-free guidance on generated texts for LVLMs.

## 3 Preliminaries

**Generative language models.** Let $p_{\boldsymbol{\theta}}$ denotes an LLM parameterized by $\boldsymbol{\theta}$. Consider a sequence $\mathbf{x} = [x_1, \ldots, x_n]$ as the input prompt, where each $x_i$ is a token from a predefined vocabulary. The LLM then generates the response sequence $\mathbf{y} = [y_1, \ldots, y_m]$ by sampling from the conditional probability distribution $p_{\boldsymbol{\theta}}(\cdot|\mathbf{x})$, where $y_t$ denotes individual token for $1 \le t \le m$. The conditional distribution

$p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{x})$ can therefore be expressed as $p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{x}) = \prod_{t=1}^{m} p_{\boldsymbol{\theta}}(y_t|\mathbf{x}, \mathbf{y}_{<t})$, where $\mathbf{y}_{<t} = [y_1, \ldots, y_{t-1}]$ for $t > 1$ and is empty for $t = 1$. In the case of LVLMs, visual tokens $\mathbf{v} = [v_1, \ldots, v_k]$ are additionally included. These tokens are generated from a pre-trained visual encoder and mapped into the token space through a linear projection. The conditional distribution of output $\mathbf{y}$ given the visual tokens $\mathbf{v}$ and textual prompt $\mathbf{x}$ is expressed as:

$$p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{v}, \mathbf{x}) = \prod_{t=1}^{m} p_{\boldsymbol{\theta}}(y_t|\mathbf{v}, \mathbf{x}, \mathbf{y}_{<t}),$$

where $p_{\boldsymbol{\theta}}$ is approximated by LVLMs.

**Guidance in generative models.** The process of a guided generation involves getting the output $\mathbf{y}$ conditioned on input $\mathbf{x}$, which encodes the desired properties of the output $\mathbf{y}$. This guidance can be generally added to the model by two distinct approaches: classifier guidance (Dhariwal and Nichol, 2021) and classifier-free guidance (Ho and Salimans, 2021). As a top-level view, both methods formulate the conditional probability distribution of output $\mathbf{y}$ conditioned on guidance $\mathbf{x}$ as

$$p(\mathbf{y}|\mathbf{x}) \propto p_{\boldsymbol{\theta}}(\mathbf{y})p(\mathbf{x}|\mathbf{y})^{\gamma}, \tag{3.1}$$

where $p_{\boldsymbol{\theta}}(\mathbf{y})$ is the original generative model and $p(\mathbf{x}|\mathbf{y})$ is the posterior distribution of $\mathbf{x}$ given $\mathbf{y}$. $\gamma$ is the guidance strength. In the classifier guidance, the posterior distribution $p(\mathbf{x}|\mathbf{y})$ in (3.1) is replaced by a classifier $p_{\boldsymbol{\phi}}(\mathbf{x}|\mathbf{y})$ parameterized by $\boldsymbol{\phi}$, which requires additional training step and calculating $\nabla_{\mathbf{x}} \log p_{\boldsymbol{\phi}}(\mathbf{x}|\mathbf{y})$. The classifier-free guidance, on the other hand, removes the necessity of the parameterized classifier $f_{\boldsymbol{\phi}}$. Instead, according to the Bayes rule, the posterior distribution can be approximated by

$$p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{y}) \propto \frac{p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{x})}{p_{\boldsymbol{\theta}}(\mathbf{y})},$$

where $p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{x})$ is the generative model when taking $\mathbf{x}$ as prompt input. Plugging (5.1) back into (3.1) yields the guided distribution that can be approximated by

$$\widehat{p}_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{x}) \propto p_{\boldsymbol{\theta}}(\mathbf{y}) \cdot \frac{p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{x})^{\gamma}}{p_{\boldsymbol{\theta}}(\mathbf{y})^{\gamma}} = \frac{p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{x})^{\gamma}}{p_{\boldsymbol{\theta}}(\mathbf{y})^{\gamma-1}}.$$

As a result, the guided LLM $\widehat{p}_{\boldsymbol{\theta}}$ places more importance on the prompt $\mathbf{x}$ during generation with the increasing value of $\gamma$, thereby producing texts that better align with the desired behavior from the prompt (Sanchez et al., 2023).

## 4 Method

To formalize the problem, we note that the architecture of LVLMs (Liu et al., 2023d; Zhu et al., 2023) is typically composed of a visual encoder, a projection layer for aligning the visual and textual domains, and an LLM for generating responses based on the image and the prompt. Object hallucination may therefore arise from deficiencies in any of these three components: (1) insufficient visual context provided by the visual encoder as highlighted by Zhang et al. (2023b) that causes a significant proportion of hallucinations, (2) flawed alignment between vision and text domains due to insufficient alignment training or the simplistic nature of the alignment layer, and (3) the inherent language priors of the LLM acquired from its pre-training data distribution (Biten et al., 2022). Aiming at the potential factors (1) and (2) that lead to object hallucinations, we introduce

5

`MARINE`, a **training-free** and **API-free** framework for mitigating LVLM hallucinations by leveraging additional object grounding features to guide the text generation of LVLMs without fine-tuning. In the following, we introduce the two major parts of `MARINE` respectively: forming the object grounding features as classify-free guidance and controlling text generation to finally mitigate object hallucinations. In Figure 2, we present the framework overview of `MARINE`.

## 4.1 Extract Object Grounding Features as Guidance

To introduce object grounding features to mitigate hallucinations, our approach integrates another object detection model DEtection TRansformer (DETR) (Carion et al., 2020), which differs from the visual encoders used in LVLM that are usually pre-trained from CLIP (Radford et al., 2021). This integration leverages DETR to extract predicted object probabilities from images, thereby providing supplementary visual information. Upon acquiring these extra visual features, we employ a "direct alignment" to the object grounding features which directly maps the output from DERT to corresponding textual objects. Direct alignment is effective and efficient. It eliminates the necessity of fine-tuning the alignment layer while retaining the full information encoded by the object grounding features. We subsequently employ a simple yet effective prompt "focusing on the visible objects in this image:" and concatenate it with the soft prompts generated from direct alignment as the classifier-free guidance prompt $\mathbf{c}$. We refrain from utilizing the hidden visual features of the DETR model but directly use the predicted object probabilities to prevent object hallucinations caused by the imperfect vision-text alignment between the DETR and LLM embedding space, as well as to eliminate the need for alignment fine-tuning.

We note that our framework is compatible with any vision model. For illustrative purposes, we utilize DETR as a representative control feature extractor and refer to our method when combined with DETR as `MARINE-DETR`. The performance of `MARINE` improves in correlation with the advancement of the control guidance extractor used. Consequently, to demonstrate the potential upper bound of `MARINE`'s performance, we consider a version utilizing a ground-truth oracle extractor, which we denote as `MARINE-Truth`.

## 4.2 Guided Text Generation

While previous classifier-free guidance method (Sanchez et al., 2023) places importance on the textual prompt itself to better align the LLM generation with user intention in the *single-modal* setting, we tackle the object hallucination problem of LVLMs by specifically placing importance on the object grounding information we introduced in the *multi-modal* setting. Therefore, in addition to the visual tokens $\mathbf{v}$ extracted from the original LVLM and textual prompt $\mathbf{x}$, we extract the auxiliary visual tokens $\mathbf{c}$ from the DETR model. The generation of the $t$-th token in the output $\mathbf{y}$ of our classifier-free guided LVLM $p_{\boldsymbol{\theta}}$ is expressed as

$$\widehat{p}_{\boldsymbol{\theta}}(y_t|\mathbf{v}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{<t}) \propto \frac{p_{\boldsymbol{\theta}}(y_t|\mathbf{v}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{<t})^{\gamma}}{p_{\boldsymbol{\theta}}(y_t|\mathbf{v}, \mathbf{x}, \mathbf{y}_{<t})^{\gamma-1}},$$

where $\mathbf{c}$ denotes our control guidance and $\gamma$ is the control strength. The sampling of output generation is given by

$$\begin{aligned}
\widehat{p}_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{v}, \mathbf{c}, \mathbf{x}) &= \textstyle\prod_{t=1}^{m} \widehat{p}_{\boldsymbol{\theta}}(y_t|\mathbf{v}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{<t}) \\
&\propto \textstyle\prod_{t=1}^{m} \frac{p_{\boldsymbol{\theta}}(y_t|\mathbf{v}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{<t})^{\gamma}}{p_{\boldsymbol{\theta}}(y_t|\mathbf{v}, \mathbf{x}, \mathbf{y}_{<t})^{\gamma-1}} \\
&= \frac{p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{v}, \mathbf{c}, \mathbf{x})^{\gamma}}{p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{v}, \mathbf{x})^{\gamma-1}}.
\end{aligned}$$

We can further view MARINE in the logit space, where the $t$-th token is therefore sampled from the logit space by

$$\log \widehat{p}_{\boldsymbol{\theta}}(y_t|\mathbf{v}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{<t}) = \gamma \log p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{v}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{<t}) + (1 - \gamma) \log p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{v}, \mathbf{x}, \mathbf{y}_{<t}).$$

This linear combination of logits implies that the conditional generation on object grounding features acts as a controllable gate. Only objects with relatively high probabilities in both branches could appear at top when sampling. Specifically, setting $\gamma = 0$ recovers the original LLM generation without control guidance and setting $\gamma = 1$ produces the LLM generation entirely based on the control. Meanwhile, for $\gamma \in (0, 1)$, MARINE yields a combination of the original generation $p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{v}, \mathbf{x})$ and the generation conditioned on the additional object grounding features $p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{v}, \mathbf{c}, \mathbf{x})$. This strikes a balance between a better ability to follow instructions to generate high-quality answers and the increased accuracy and detail in image descriptions. We summarize MARINE in Algorithm 1.

---

**Algorithm 1** MARINE

---

1: **Input:** LLM with parameter $\boldsymbol{\theta}$, input prompt $\mathbf{x}$, visual tokens $\mathbf{v}$ from LVLM's original vision tower, auxiliary visual tokens $\mathbf{c}$ from the DETR model,
2: Initialize empty output $\mathbf{y} = []$.
3: **for** $t = 0, 1, ...$ **do**
4:      Construct unconditional input $\mathbf{x}_{\text{uncond}}^{(t)} = [\mathbf{v}, \mathbf{x}, \mathbf{y}_{<t}]$.
5:      Generate unconditional output logits using LLM: $\ell_{\text{uncond}}^{(t)} = \log p_{\boldsymbol{\theta}}(\mathbf{x}_{\text{uncond}}^{(t)})$.
6:      Construct conditional input $\mathbf{x}_{\text{cond}}^{(t)} = [\mathbf{v}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{<t}]$.
7:      Generate conditional output logits using LLM: $\ell_{\text{cond}}^{(t)} = \log p_{\boldsymbol{\theta}}(\mathbf{x}_{\text{cond}}^{(t)})$.
8:      Update output logits $\ell^{(t)} = \gamma \ell_{\text{cond}}^{(t)} + (1 - \gamma)\ell_{\text{uncond}}^{(t)}$.
9:      Sample token $y_t$ from logit space denoted by $\ell^{(t)}$.
10:     Let $\mathbf{y} = [\mathbf{y}, y_t]$.
11:     **if** $y_t = $ <EOS> **then**
12:        **break**
13:     **end if**
14: **end for**
15: **Output:** $\mathbf{y}$.

---

# 5   Experiments

In this section, we provide detailed empirical evaluations of our method MARINE in mitigating object hallucinations across a variety of LVLMs. Specifically, we highlight that MARINE outperforms the state-of-the-art mitigation methods on established metrics across different question formats.

## 5.1   Experiment Setup

**Models.** To demonstrate the broad applicability of our approach across different LVLM architectures, we apply and evaluate MARINE to recent widely-used models including *LLaVA* (Liu et al., 2023d), *LLaVA-v1.5* (Liu et al., 2023c), *MiniGPT-v2* (Chen et al., 2023), *mPLUG-Owl2* (Ye et al., 2023), *InstructBLIP* (Liu et al., 2023c) and *LLaMA-Adapter-v2* (Gao et al., 2023). To address the object hallucination problems in text generation, we incorporate the DEtection TRansformer (DETR) (Carion et al., 2020) as an object grounding encoder to enrich the visual features.

**Datasets and evaluations.** In alignment with established evaluations from previous studies (Dai et al., 2023b; Yin et al., 2023), we assess our method using the following metrics:

- Caption Hallucination Assessment with Image Relevance (*CHAIR*) (Rohrbach et al., 2018). It involves prompting the LVLMs to generate a description for the input image, and then comparing this generation with ground truth objects present in the image. CHAIR quantifies hallucination both at instance level and sentence level, respectively defined as $\text{CHAIR}_I$ and $\text{CHAIR}_S$:

$$\text{CHAIR}_I = \frac{\left|\{\text{hallucinated objects}\}\right|}{\left|\{\text{all mentioned objects}\}\right|}, \quad \text{CHAIR}_S = \frac{\left|\{\text{captions with hallucinated objects}\}\right|}{\left|\{\text{all captions}\}\right|}.$$

  In addition to these metrics, we incorporate an instance-level Recall score in our evaluation to evaluate whether the descriptions accurately include the necessary visual content from the image:

$$\text{Recall} = \frac{\left|\{\text{non-hallucinated objects}\}\right|}{\left|\{\text{all existing objects}\}\right|}.$$

- Polling-based Object Probing Evaluation (*POPE*) (Li et al., 2023b). POPE formulates a binary classification task by prompting LVLMs with questions such as "Is there a keyboard in this image?" to answer "yes" or "no". We specifically chose the adversarial setting, which is considered the most challenging setting. We report the accuracy and F1 score of the LVLMs' responses, and the proportion of "yes" answers in this experiment.
- *GPT-4V-aided Evaluation* (Yin et al., 2023). The GPT-4V-aided evaluation compares the outputs of two LVLM assistants using GPT-4V as a judge. GPT-4V is asked to provide scores out of 10 on the following two metrics: 1) *accuracy*: how accurately each assistant describes the image; and 2) *detailedness*: the richness of necessary details in the response. In this evaluation, we utilize the LLaVA-QA90 task (Liu et al., 2023d)[2] and additionally consider the image captioning task.

All evaluations are conducted using the MSCOCO val2014 dataset. Consistent with Li et al. (2023b), we use the same random subset of 500 images for both the CHAIR and POPE evaluations. For the GPT-4V-aided evaluation, we additionally use 90 questions from the LLaVA-QA90 task and additionally selected a sample of 50 images for the image captioning task on MSCOCO val2014 dataset.

**Baselines.** In addition to comparing with the performance of the original LVLM sampling method, we also consider the following popular methods for mitigating hallucinations.

- *Greedy-Decoding*, which adopts the greedy sampling strategy, by generating tokens with the highest posterior probability to address hallucinations arising from.
- LVLM Hallucination Revisor (*LURE*) (Zhou et al., 2023), which identifies and masks potentially hallucinated words and fine-tune a MiniGPT4 model to rectify object hallucinations in the generated descriptions.
- *LURE with Cutoff*. As demonstrated in Figure 1, the original LURE method tends to generate long descriptions regardless of the provided instructions, which sometimes results in worse performance of CHAIR as unnecessary information is included. Therefore, we also introduce a modified baseline, where we truncate the LURE's output to match the length (in terms of the number of sentences) of the original generations.
- Visual Contrastive Decoding (*VCD*), which distorts the image inputs to impose penalties on logit outputs. As VCD only evaluated with POPE on LLaVA-v1.5 and InstructBLIP, we compare with their reported performances.

---

[2] https://github.com/haotian-liu/LLaVA/blob/main/playground/data/coco2014_val_gpt4_qa_30x3.jsonl

- *Woodpecker* (Yin et al., 2023), which leverages GPT-3.5 to correct hallucinations in LVLM generation with five steps toward the correction.

Further experiment details on model architectures, datasets and evaluations are deferred to Appendix A.

## 5.2 Results

Our findings are comprehensively presented in Table 1 and Table 2, where we compare `MARINE` with the baselines on the CHAIR and POPE metrics. We highlight that, while previous mitigation methods focus on the initial versions of the LVLMs (e.g., LLaVA), our experiments further encompass their latest versions (e.g., LLaVA-v1.5) which have been better trained for fewer hallucinations. Overall, `MARINE` achieves superior performances across different LVLM architectures and evaluation metrics, ranking as the best or second-best on the majority of the tasks.

In Table 1, we present the CHAIR evaluation, where `MARINE` achieves a substantial improvement up to +22.0% on CHAIR$_S$ and +35.2% on CHAIR$_I$ compared to the original outputs. Notably, while previous baselines (specifically the fine-tuned LURE and GPT-3.5-aided Woodpecker) are ineffective in reducing hallucinations on the latest LVLM versions, which already exhibit decent performances on CHAIR, `MARINE` provides even further improvements. On the mPLUG-Owl2 model, for example, CHAIR$_S$ and CHAIR$_I$ scores have markedly improved from 5.9 and 3.5 to 2.8 and 1.2, respectively. Furthermore, `MARINE` significantly enhances the LVLMs' ability to focus on existing objects, leading to more detailed responses, as evidenced by the marked average increase of +40% on Recall. When utilizing `MARINE`, LLaVA achieves a notable recall score of 67.8%.

The POPE evaluation, detailed in Table 2, further validates the superior performance of `MARINE` against existing baselines on different question formats. `MARINE` consistently outperforms the original outputs by a large margin, achieving improvements of up to +21.4% on accuracy and +12.0% on F1 score. Specifically, for LLaVA model, `MARINE` realizes gains of +28.8% on accuracy and +15.0% on F1 score. Moreover, `MARINE` significantly outperforms VCD and Woodpecker, with the margins exceeding +3.8%. It also significantly reduces the biased responses in LVLMs toward answering "yes", as evidenced by the closer-to-50% "yes" ratio (a 22.4% shift in average towards non-biased answers). This indicates that `MARINE` not only mitigates hallucinations but also, to a notable extent, addresses the common issue of overconfidence which often leads to a "yes" bias in LVLM responses.

Table 1: Evaluation with CHAIR score across multiple LVLM architectures comparing our method with several baselines. We report CHAIR$_I$, CHAIR$_S$ and the recall score. The **bold** numbers indicate the best results among the methods evaluated and the underscored numbers represent the second-best results. We show `MARINE-Truth` as an ideal reference performance of `MARINE`.

| Method | LLaVA | | | LLaVA-v1.5 | | | MiniGPTv2 | | | mPLUG-Owl2 | | | InstructBLIP | | | LLaMA-Adapter-v2 | | |
| CHAIR | $C_S\downarrow$ | $C_I\downarrow$ | $R\uparrow$ | $C_S\downarrow$ | $C_I\downarrow$ | $R\uparrow$ | $C_S\downarrow$ | $C_I\downarrow$ | $R\uparrow$ | $C_S\downarrow$ | $C_I\downarrow$ | $R\uparrow$ | $C_S\downarrow$ | $C_I\downarrow$ | $R\uparrow$ | $C_S\downarrow$ | $C_I\downarrow$ | $R\uparrow$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Original | 26.1 | 10.8 | 46.1 | **7.5** | <u>4.2</u> | 40.9 | 11.6 | 6.8 | 36.5 | <u>5.9</u> | 3.5 | 36.1 | 5.3 | 3.6 | 31.2 | 27.1 | 10.5 | 50.9 |
| Greedy | 26.6 | 10.5 | 47.4 | 8.8 | 4.6 | 41.1 | <u>8.2</u> | **4.2** | 41.1 | 6.2 | <u>3.4</u> | 38.8 | <u>3.6</u> | <u>2.3</u> | 34.3 | 26.6 | 10.4 | 50.8 |
| LURE | 33.8 | 11.6 | <u>54.8</u> | 38.9 | 11.2 | <u>56.3</u> | 36.2 | 11.4 | <u>54.6</u> | 33.9 | 10.8 | **55.9** | 38.1 | 12.1 | **54.5** | 35.2 | 10.7 | <u>55.2</u> |
| LURE w/ cutoff | 24.4 | 9.3 | 50.2 | 18.4 | 6.8 | 47.3 | 12.5 | 6.2 | 42.0 | 15.4 | 6.6 | 45.5 | 9.6 | 6.4 | 34.5 | 27.8 | 9.5 | 51.3 |
| Woodpecker | **19.5** | <u>8.9</u> | 44.3 | <u>8.5</u> | 4.5 | 38.4 | **7.5** | <u>4.5</u> | 37.0 | 8.0 | 4.3 | 37.5 | 8.0 | 6.2 | 32.6 | <u>22.0</u> | <u>7.8</u> | 49.0 |
| `MARINE-DETR` | <u>24.3</u> | **7.1** | **67.8** | 7.5 | **3.7** | **60.7** | 10.0 | 4.7 | **62.2** | **4.8** | **2.4** | <u>50.0</u> | **2.8** | **1.6** | <u>37.8</u> | **14.8** | 5.5 | **59.3** |
| `MARINE-Truth` | 13.8 | 3.1 | 92.3 | 6.4 | 2.1 | 71.0 | 8.9 | 2.4 | 76.2 | 2.8 | 1.2 | 54.4 | 2.3 | 1.3 | 34.4 | 13.2 | 6.6 | 53.8 |

**Discussion on fine-tuning methods.** The examples depicted in Figure 1 illustrate that LURE, at times, fails to adhere to the given instructions when correcting LVLM generations. Despite receiving

Table 2: Evaluation with POPE score in adversarial setting across multiple LVLM architectures comparing our method with several baselines. We report the POPE accuracy (%), F1 score (%) and the yes ratio (%). The ideal yes ratio for a non-biased LVLM is 50%. The **bold** numbers indicate the best results among the methods evaluated and the underscored numbers represent the second-best results. We show `MARINE-Truth` as an ideal reference performance of `MARINE`.

| Method POPE | LLaVA Acc ↑ | F1 ↑ | Yes | LLaVA-v1.5 Acc ↑ | F1 ↑ | Yes | MiniGPTv2 Acc ↑ | F1 ↑ | Yes | mPLUG-Owl2 Acc ↑ | F1 ↑ | Yes | InstructBLIP Acc ↑ | F1 ↑ | Yes | LLaMA-Adapter-v2 Acc ↑ | F1 ↑ | Yes |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Original | 53.5 | 68.1 | 95.7 | 79.0 | 81.1 | 61.2 | 80.2 | 79.0 | 44.5 | 71.5 | 76.6 | 71.9 | 71.6 | 74.7 | 61.4 | 56.0 | 68.9 | 91.4 |
| Greedy | 51.8 | 67.4 | 97.7 | 79.4 | _81.6_ | 61.6 | _82.7_ | _81.7_ | 44.5 | 72.5 | _77.5_ | 72.4 | _79.8_ | _81.4_ | 58.6 | 55.7 | 68.8 | 92.1 |
| LURE | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
| VCD (reported) | - | - | - | _80.9_ | 81.3 | - | - | - | - | - | - | - | 79.6 | 79.5 | - | - | - | - |
| Woodpecker | _77.5_ | _77.6_ | **50.5** | 80.5 | 80.6 | **50.5** | 79.5 | 77.8 | 42.5 | _77.5_ | 76.9 | **47.5** | 79.0 | 78.6 | _48.0_ | 77.0 | 77.2 | **51.0** |
| `MARINE-DETR` | **80.6** | **82.4** | _60.1_ | **90.2** | **90.3** | _51.3_ | **87.7** | **88.1** | 53.5 | **77.9** | **81.4** | _69.1_ | **90.7** | **90.5** | 48.2 | _66.6_ | _69.4_ | _59.2_ |
| `MARINE-Truth` | 84.9 | 86.4 | 61.1 | 96.1 | 96.2 | 53.5 | 90.8 | 91.2 | 54.0 | 80.0 | 83.4 | 69.9 | 96.9 | 96.9 | 50.2 | 68.6 | 71.4 | 59.7 |

concise image descriptions generated based on instructions for short responses, LURE predominantly overwrites them with excessively long responses that contain information irrelevant to the instruction. Furthermore, LURE fails to adequately address the binary question format of POPE, as LURE fixates on extended descriptions without responding with "yes" or "no", making its evaluation using POPE impractical. This issue can be prevalent in small-scale fine-tuning methods, where the limited variety of the specifically tailored fine-tuning dataset harms the model's performance on other tasks. In contrast, the training-free approach of `MARINE` demonstrates effective mitigation of hallucinations across a variety of question formats.

**Results on GPT-4V-aided evaluation.** In addition to CHAIR and POPE, we assess `MARINE` on LLaVA-QA90, a comprehensive visual question-answering task including conversations, visual perceptions, and complex reasoning. Following Yin et al. (2023), we utilize GPT-4V to evaluate and compare the performance of the original LVLMs and LVLMs with `MARINE`. To formulate the evaluation prompt for GPT-4V[3], we provide the original question as well for a task-orientated evaluation. The prompt used for assessment is in Appendix A. As shown in Table 3, `MARINE` achieves superior performance gains on GPT-4V-aided evaluation in both accuracy and detailness metrics. In addition, we leverage GPT-4V to assess `MARINE` on the image captioning task of MSCOCO, where we observe consistent enhancements on this task.

## 5.3 Ablation Study

In this study, we explore the effect of guidance strength and the impact of noise intensity of object grounding features on mitigating object hallucinations in LVLMs through both quantitative and qualitative analysis. Additionally, we present concrete examples to demonstrate the influence of control guidance on the output logits of the LVLMs.

**Effect of guidance strength.** In Figure 3 and Figure 4, we demonstrate the influence of the guidance strength on the CHAIR metrics for LLaVA and InstructBLIP models, focusing on the effectiveness of reducing object hallucinations. An increase in guidance strength from 0 to 1 leads to a notable decrease in CHAIR scores, particularly in $CHAIR_I$. This trend implies that increasing guidance strength significantly reduces hallucinated outputs, with an average reduction of 22.0% on $CHAIR_S$ and 35.2% on $CHAIR_I$, thereby enhancing the models' ability to produce accurate

---

[3]We used `gpt-4-1106-vision-preview` in obtaining our final experiment results. As OpenAI continues to update its API, different versions may result in slightly different values.

Table 3: Results of GPT-4V-aided evaluation. The accuracy and detailedness metrics are on a scale of 10, and a higher score indicates better performance.

| Method | w/MARINE | LLaVA-QA90 | | Image Captioning | |
|---|---|---|---|---|---|
| | | Acc ↑ | Detail ↑ | Acc ↑ | Detail ↑ |
| LLaVA | ✗ | 5.6 | 4.4 | 5.2 | 4.6 |
| | ✓ | **6.0** | **4.7** | **5.6** | **4.7** |
| LLaVA-v1.5 | ✗ | 6.9 | 5.1 | 7.4 | 5.8 |
| | ✓ | **7.3** | **5.4** | **7.7** | **6.2** |
| MiniGPTv2 | ✗ | 6.6 | 3.1 | 6.2 | 4.9 |
| | ✓ | **6.6** | **4.2** | **6.9** | **5.5** |
| mPLUG-Owl2 | ✗ | 6.3 | 4.8 | 7.0 | 5.5 |
| | ✓ | **6.8** | **5.0** | **8.1** | **6.2** |
| InstrcutBLIP | ✗ | 5.9 | 4.8 | 7.0 | 5.5 |
| | ✓ | **6.5** | **5.2** | **7.8** | **6.0** |
| LLaMA-Adapter-v2 | ✗ | **6.1** | **4.8** | 6.0 | 5.0 |
| | ✓ | 5.0 | 3.2 | **6.7** | **5.2** |

descriptions. It's crucial to note that, although some models exhibit optimal performance at a guidance strength of $\gamma = 1$, excessively strong guidance can adversely affect the models' ability to adhere to provided instructions, as evidenced in Figure 5. This observation highlights the necessity of having a balanced guidance strength that ensures high-quality, accurate outputs while adhering closely to the given instructions. Based on our findings, we recommend a guidance strength within the range of $\gamma \in (0.3, 0.7)$ as the most effective for maintaining this balance.



(a) CHAIR$_S$  (b) CHAIR$_I$  (c) Recall

Figure 3: LLaVA's performance on CHAIR according to different guidance strength $\gamma$ of MARINE.

**Impact of noise intensity of object grounding features.** In Figure 6, we delve into the impact of the quality of the object grounding features in MARINE on the performance of LVLMs. We maintain a constant guidance strength of 0.5 and 1.0 while varying the object grounding features at five distinct levels of noise intensity. This variation is achieved by implementing four confidence thresholds (0.5, 0.7, 0.9, and 0.95) in the DETR model predictions, where higher thresholds correspond to lesser, yet higher-quality, visual information. As a comparative standard, we include MARINE-Truth as an ideal reference performance of MARINE. Our findings highlight two significant insights. Firstly, an increase
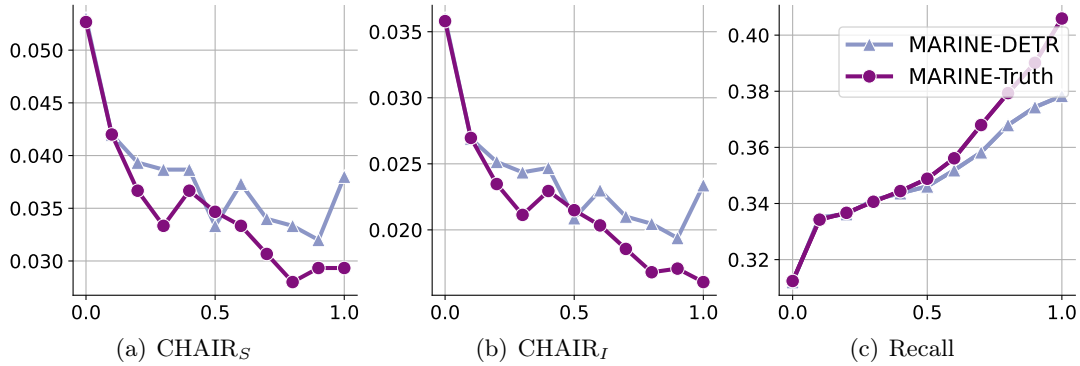
(a) CHAIR$_S$  (b) CHAIR$_I$  (c) Recall

Figure 4: InstructBLIP's performance on CHAIR according to different guidance strength $\gamma$ of `MARINE`.



What precautions has the little girl taken to stay dry in the rain?

**LLaVA w/ MARINE γ = 0.5**
The little girl is holding an umbrella to protect herself from the rain, ensuring that she stays dry during the rainy weather.

**LLaVA w/ MARINE γ = 1.0**
The little girl has taken a red and white umbrella out of her backpack and is holding it up to protect herself from the rain. She is standing in front of a truck, possibly waiting for her mother or guardian to pick her up. The umbrella provides a visible and practical means of protection against the rainfall, allowing the girl to stay dry and comfortable while waiting for her ride.

Figure 5: An example of the negative impact of excessive guidance on LVLM's ability to follow instructions accurately. While the response with $\gamma = 1$ identifies more existing objects, it introduces irrelevant information to the instruction.

in the quality of visual information correlates with a noticeable decrease in hallucinations produced by the LVLMs. A lower threshold, which allows for more visual information but also includes noisier content, could potentially result in an increased occurrence of hallucinations. Furthermore, while lower-quality visual information may not effectively reduce hallucinations, it is associated with enhanced Recall. This suggests that LVLMs under guidance, despite the presence of noisy visual inputs, tend to focus more on the visual details (i.e., objects), resulting in more elaborate descriptions.

**Effect of `MARINE` on logit distribution.** In Figure 7, we illustrate a specific example that shows how `MARINE` influences the logit distribution of LVLMs during text generation. Specifically, `MARINE` is observed to selectively target the potential hallucinated tokens, reducing their original probabilities to mitigate the risk of hallucination in the generated text. For instance, in the provided example, the

Figure 6: LLaVA's performance on CHAIR according to different noise intensity of object grounding features in `MARINE`. We consider four confidence thresholds (0.5, 0.7, 0.9, and 0.95) for DETR to vary the noise intensity.

probability of "fork" is significantly lowered with `MARINE`, which would have originally resulted in a hallucinated object. Conversely, standard language elements such as "various", an adjective describing the overall image context, and "with", a crucial preposition, maintain their original probabilities. This selective nature of modulation by `MARINE` ensures coherent and contextually relevant text generation that adheres to the instruction while effectively reducing hallucinations.

# 6 Conclusion and Future Work

In this paper, we introduced a training-free and API-free framework `MARINE` to mitigate object hallucination in LVLMs during its text generation process. Leveraging a pre-trained object grounding vision encoder for a novel classifier-free guidance framework in the multi-modal setting, `MARINE` effectively and cost-efficiently reduces the hallucinations of six widely-used LVLMs, as assessed by various metrics across different tasks. The inherent compatibility of the `MARINE` with various vision models and projection functions further underscores its flexibility. In contrast to post-generation correction methods, `MARINE` strikes a balance between efficiency, instruction-following ability and effectiveness in reducing object hallucinations.

**Limitations and future work.** `MARINE` exhibited impressive performance with the DETR object grounding encoder. However, there is potential for enhancement by incorporating more advanced vision encoders and investigating their varying impacts within our framework. Additionally, although extensive experiments have been conducted on various LVLM architectures and evaluation metrics, further evaluation of `MARINE` across a broader range of benchmarks would be advantageous.
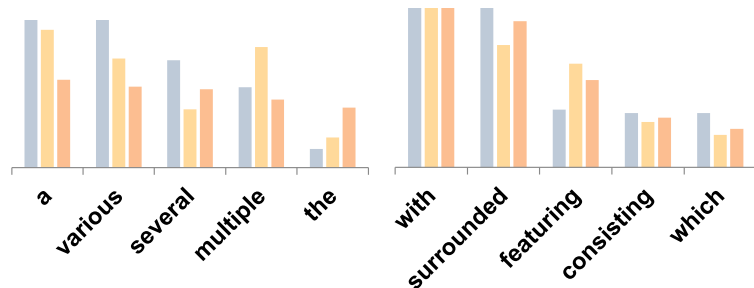
# A Experiment Details

## A.1 Model Architectures

In Table 4, we provide detailed descriptions of the LVLM architectures used in our experiments. These LVLMs respectively leverage the pre-trained vision encoder of the models we listed, which are all based on the Vision Transformer (ViT) (Dosovitskiy et al., 2020) architecture.

13

(a) An example of image description where the original LLaVA outputs a hallucinated object, "fork".



(b) The probability distributions at the token of the hallucinated word in the original, control, and MARINE outputs. MARINE effectively decrease the the probability of "fork".



(c) Probabilities of non-hallucinated words remain the same, highlighting MARINE's ability to preserve normal outputs.

Figure 7: This sample shows how MARINE controls logit distributions to mitigate hallucinations like "fork" while preserving the probabilities of "with", "various" during generation.

Table 4: Details of the LVLM architectures that we used in our paper.

| Model | Vision encoder | LLM |
|---|---|---|
| LLaVA (Liu et al., 2023d) | CLIP-L (Radford et al., 2021) | LLaMA-2-7B-Chat (Touvron et al., 2023b) |
| LLaVA-v1.5 (Liu et al., 2023c) | CLIP-L-336px (Radford et al., 2021) | Vicuna-v1.5-7B (Chiang et al., 2023) |
| MiniGPT-v2 (Chen et al., 2023) | EVA-G (Fang et al., 2023) | LLaMA-2-7B-Chat (Touvron et al., 2023b) |
| mPLUG-OWL2 (Ye et al., 2023) | CLIP-L (Radford et al., 2021) | LLaMA-2-7B (Touvron et al., 2023b) |
| InstructBLIP (Dai et al., 2023a) | BLIP-2 (Li et al., 2023a) | Vicuna-v1.1-7B (Chiang et al., 2023) |
| LLaMA-Adapter-v2 (Gao et al., 2023) | CLIP-L-336px (Radford et al., 2021) | LLaMA-7B (Touvron et al., 2023a) |

## A.2 Experiment Setting for Hallucination Evaluations

Key factors that potentially affect the hallucination evaluation outcomes, including the evaluation dataset and prompt template, LVLM's sampling strategy and batched generation techniques, and guidance strength, are detailed in this section.

**Experiment setting for CHAIR evaluation.** We adopt the same prompt "Generate a short caption of the image." as utilized by Li et al. (2023b). Across the six LVLMs tested, we standardized the sampling strategy with temperature set to 0.6 and top-p to 0.9. Each test was conducted multiple times using three distinct random seeds (42, 142, 242), with the mean score reported in our results. Additionally, we employed the batched generation to expedite the evaluation process. We avoid the negative impact of batched generation by adopting left padding as it leads to more accurate results if the LVLM does not explicitly assign the padding strategy for inference. For the calculation of CHAIR metrics, we referenced the 80 object categories annotated in the MSCOCO dataset, following Rohrbach et al. (2018). Besides, we employed the synonym list from Lu et al. (2018) to align synonymous words in the generated text with MSCOCO object categories. Additionally, due to the cost considerations associated with the GPT-3.5 API, we limited our analysis to 200 samples for Woodpecker correction for each model and reported the result in Table 1.

**Experiment setting for POPE evaluation.** POPE is a flexible approach to evaluating hallucinations in LVLMs, which formulates a binary classification task by prompting LVLMs with questions such as "Is there a keyboard in this image?" to answer "yes" or "no". We specifically chose the adversarial settings, the most challenging setting, which constructs POPE questions from the top-k most frequently co-occurring but absent objects. Following Li et al. (2023b), we used the same 500 images from the MSCOCO val2014 dataset, applying 6 POPE questions per image for POPE evaluation. Similarly, we constrained our analysis to 200 samples for Woodpecker correction for each model due to the high costs associated with the GPT API. The outcomes of this analysis are detailed in Table 2.

**Experiment setting for GPT-4V-aided evaluation.** As shown in Figure 8, the assessment prompt template we used is slightly different from that of Yin et al. (2023). Specifically, we also include the original question for a task-orientated evaluation and exclude prompts that describe Woodpecker-specific output formats like object bounding boxes. Examples of the GPT-4V-aid evaluation responses are illustrated in Figure 9 and 10. Besides, a fixed guidance strength of 0.5 was used in the evaluations in Table 3. Utilizing the `gpt-4-1106-vision-preview`, all final experiments were conducted between 01/01/2024-01/30/2024. As OpenAI continues to update its API, accessing different versions may result in slightly different values.

**Experiment setting for ablation study.** To quantitatively assess the influence of guidance strength, we varied it from 0 to 1, as shown in Figure 7. Additionally, to investigate the impact of the noise intensity of object grounding features, we conducted experiments on five levels of noise intensity as depicted in Figure 6. These quantitative experiments were conducted using the same setting as those in CHAIR evaluation. For qualitative analysis, exemplified in Figure 5 and 7, we selected guidance strength from a recommended range of $\gamma \in (0.3, 0.7)$.

## A.3 Additional Experiment Results

Additional experimental results explore the impact of guidance strength (ranging from 0 to 1) and the noise intensity of object grounding features (across four confidence thresholds: 0.5, 0.7, 0.9, and 0.95, with `MARINE-Truth` serving as an ideal reference). These effects are examined across LLaVA,

Figure 8: Prompt template for GPT-4V-aided evaluation. {question} is the original instruction; {answer 1} is the original response, and {answer 2} is the response generated by the LVLM using `MARINE-DETR` with a guidance strength of 0.5.

InstructBLIP, and mPLUG-Owl2, with findings presented in Figures 11, 12, 13, 14, and 15.

## A.4   Generation Examples

In Figures 9 and 10, we present examples of GPT-4V-aided evaluations based on the outputs of LLaVA-v1.5 and LLaVA-v1.5 with `MARINE`. In Figures 16, 17, 18, 19, and 20, we present examples of the outputs from LURE (Zhou et al., 2023), Woodpecker (Yin et al., 2023) and `MARINE` on different tasks further validate our arguments in the paper.

**Prompt**
Question:
Can you elaborate on the elements of the picture provided?

Assistant 1 (LLaVA-v1.5):
The image features a desk with various items on it. There is a bookshelf with a few books, including a prominent one titled \"Administrative Law \" A blue and white tie is placed on the desk, adding a touch of color and formality to the scene. A pair of sunglasses is also resting on the desk, possibly belonging to someone who was recently there.

In addition to the books and tie, there is a cup and an apple on the desk, indicating that the person might have been enjoying a snack or drink while working or studying. A teddy bear is also present in the scene sitting on the desk, and adding a playful touch to the setting.

Assistant 2 (LLaVA-v1.5 w/ MARINE):
The image features a desk with various items on it. There is a book, possibly a law book, positioned on the left side of the desk. A blue and white tie is placed next to the book, adding a touch of color to the scene. On the right side of the desk, there is an apple, possibly a Granny Smith apple, and a pair of sunglasses.

Additionally, there is a stuffed teddy bear located at the top right corner of the desk, giving the scene a cozy and personal touch. A baseball glove is also present in the image, suggesting that the owner might have an interest in sports.

**GPT-4V Evaluation**
Accuracy:
Scores of the two answers: 4 7
Reason: Both assistants mention items that are present in the image, such as the book with the title "Administrative Law," the tie, the apple, the sunglasses, and the teddy bear. Assistant 1 hallucinates the presence of a cup, which is not visible in the image. Assistant 2 inaccurately describes the presence of a baseball glove, which is not in the image; the item mistaken for a glove is actually part of the teddy bear's outfit. Assistant 2 also avoids mentioning the specific title of the book, even though it is visible in the image.

Detailedness:
Scores of the two answers: 5 6
Reason: Assistant 1 provides a more detailed description of the scene by elaborating on the potential reasons for the items' presence, such as the person enjoying a snack or the teddy bear adding a playful touch, However, the mention of a cup, which is not in the image, detracts from the quality of the response. Assistant 2 provides a sufficient amount of detail without over-elaborating and makes a specific guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score.

Figure 9: An example of GPT-4V-aided evaluation on a detailed question.

**Prompt**

Question:

What might be the purpose of this table arrangement?

Assistant 1 (LLaVA-v1.5):

The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation and snacking. The presence of the book and sunglasses suggests that the person might be taking a break from reading or working on a project, while the apple provides a healthy snack option. The combination of these items on the table indicates that the person values both work and relaxation, aiming to create a balance between the two.

Assistant 2 (LLaVA-v1.5 w/ MARINE):

The purpose of this table arrangement appears to be a combination of work and leisure. The presence of books. a sunglasses. a tie. and an apple suggests that the person using this space may be engaged in work or study, as indicated by the books. Additionally, the apple and sunglasses indicate that the person might be taking a break, enjoying a snack, or preparing for a more relaxed activity. The presence of the teddy bear adds a touch of playfulness and personalization to the setting, creating a cozy and inviting atmosphere.

**GPT-4V Evaluation**

Accuracy:

Scores of the two answers: 6 7

Reason: Both Assistants made a few errors in their descriptions, Assistant 1 incorrectly states that there\'s a book on the table, which is ambiguous; technically, the books are not on the table but more likely on a shelf next to it and the tie isn't mentioned, which is a significant element in the image. Assistant 2 correctly recognizes the presence of books, a tie, sunglasses, and an apple, along with the teddy bear dressed in a sports jersey. However Assistant 2 also refers to books in a slightly ambiguous way, implying they are on the table, and incorrectly describes the sunglasses as "a sunglasses" which is grammatically incorrect and could slightly mislead the context.

Detailedness:

Scores of the two answers: 6 7

Reason: Assistant 1 provides an explanation of the potential purpose of the table arrangement but misses out mentioning the tie and teddy bear both of which add context to the scenario depicted. The description lacks detail on the actual arrangement and how it specifically suggests a mixture of work and leisure. Assistant 2 includes more elements in the description, such as the tie and teddy bear, adding to the detail and mentioning relaxation and personalization, which gives a slightly better insight into the purpose of the table arrangement, but their response could also be more nuanced by acknowledging the role of each item.

Figure 10: An example of GPT-4V-aided evaluation on a complex question.
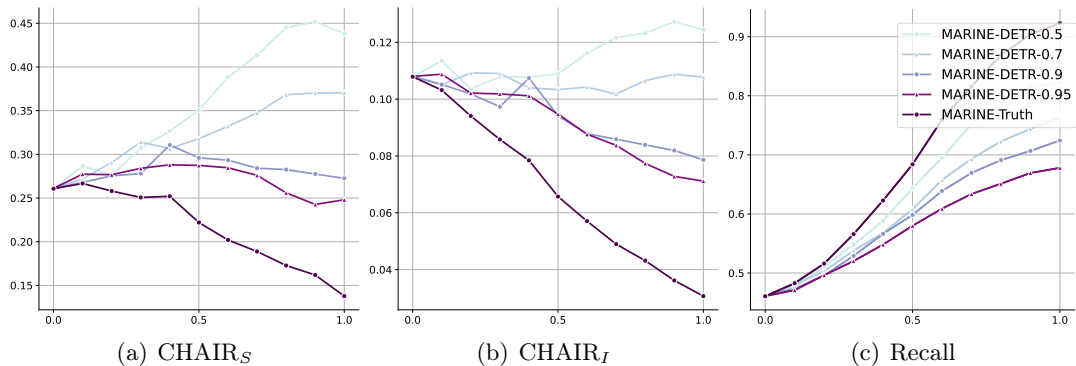


(a) CHAIR$_S$  (b) CHAIR$_I$  (c) Recall

Figure 11: LLaVA's performance on CHAIR according to different guidance strength and noise intensity of object grounding features in MARINE.
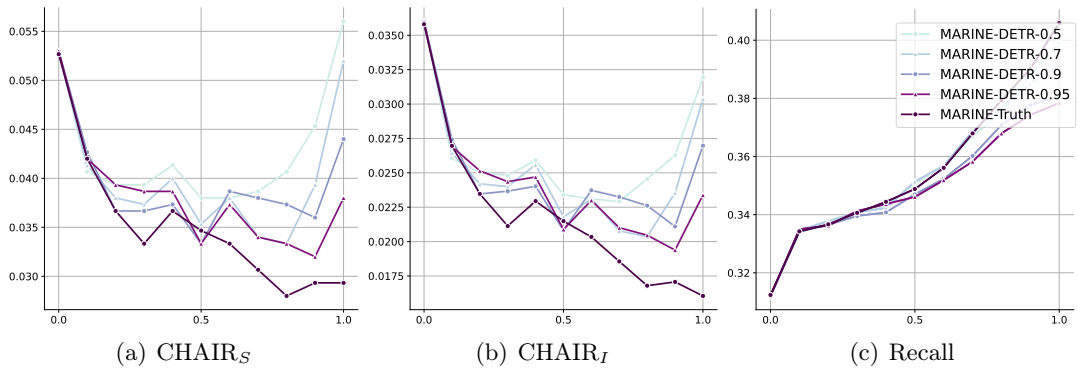
Figure 12: InstrucBLIP's performance on CHAIR according to different guidance strength and noise intensity of object grounding features in `MARINE`.
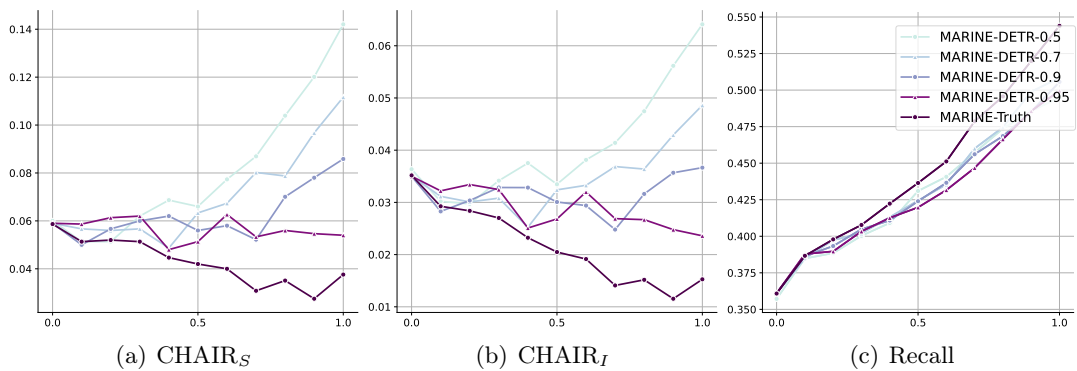


Figure 13: mPLUG-Owl2's performance on CHAIR according to different guidance strength and noise intensity of object grounding features in `MARINE`.
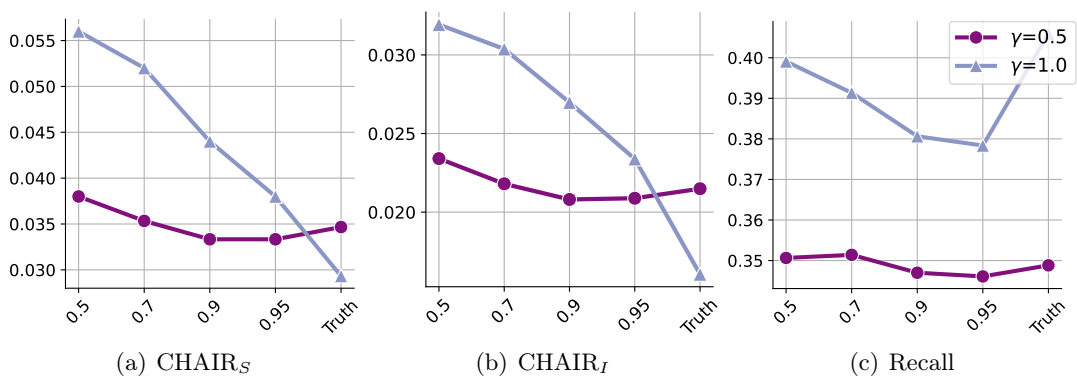


Figure 14: InstructBLIP's performance on CHAIR according to different noise intensity of object grounding features in `MARINE`. We consider four confidence thresholds (0.5, 0.7, 0.9, and 0.95) for DETR to vary the noise intensity, with `MARINE-Truth` serving as an ideal reference.

| (a) CHAIR$_S$ | (b) CHAIR$_I$ | (c) Recall |

Figure 15: mPLUG-Owl2's performance on CHAIR according to different noise intensity of object grounding features in `MARINE`. We consider four confidence thresholds (0.5, 0.7, 0.9, and 0.95) for DETR to vary the noise intensity, with `MARINE-Truth` serving as an ideal reference.
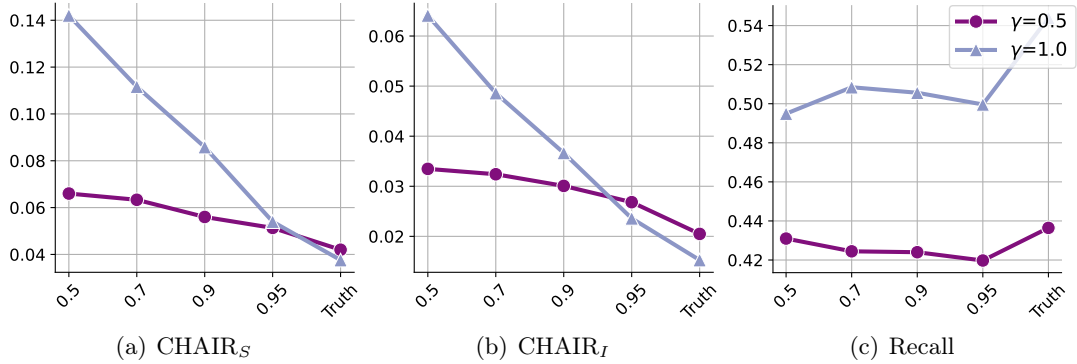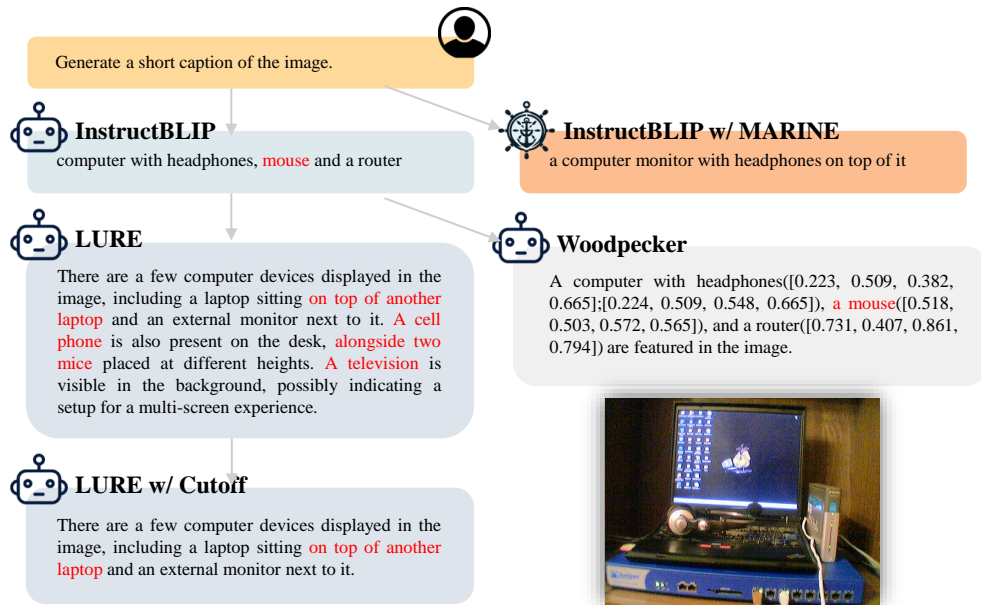


Figure 16: A comparison of responses from baseline models and our `MARINE` in an image description task. It illustrates `MARINE`'s superior ability to reduce hallucinations, in contrast to LURE and Woodpecker, which fail to effectively address hallucinations and sometimes even increase hallucinated content. This example highlights the strengths of our correct-during-generation framework over post-correction approaches, showcasing its efficiency, preservation of original style, and enhanced adherence to instructions.
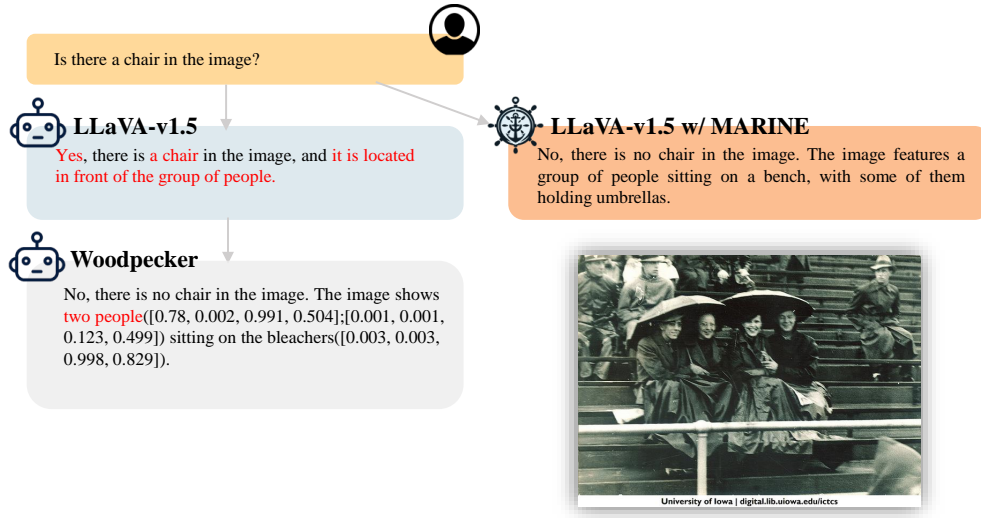
Figure 17: A comparison of responses from Woodpecker and our MARINE in POPE "yes-or-no" task.
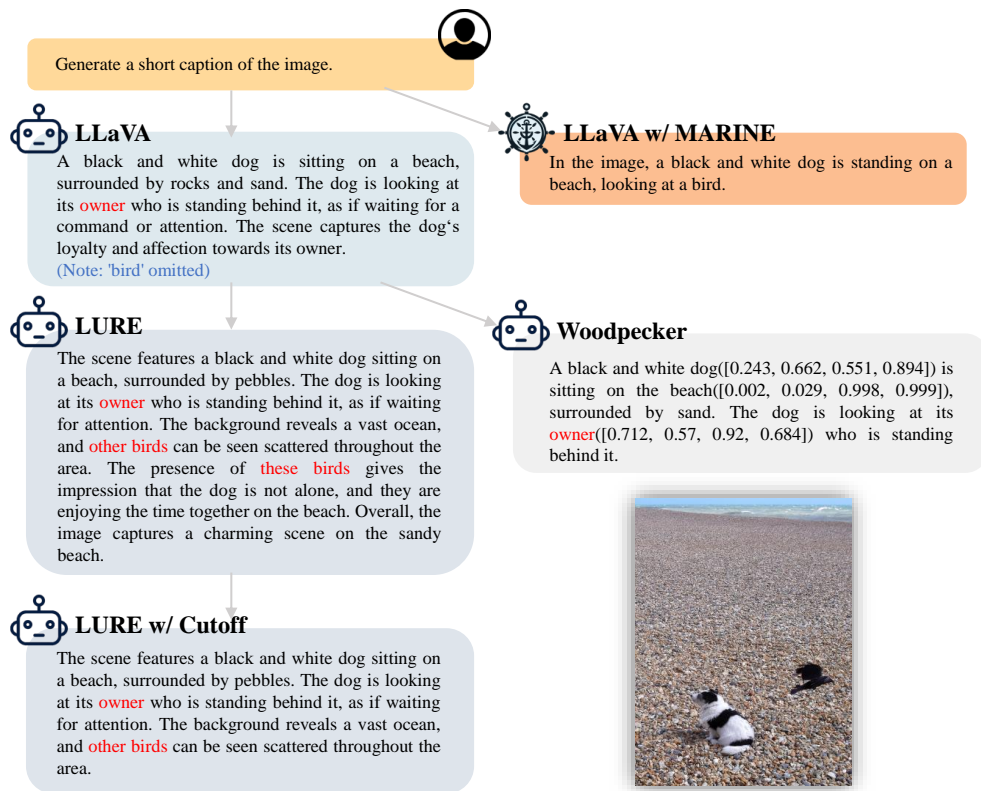


Figure 18: A comparison of responses from baseline models and our MARINE in an image description task. MARINE effectively reduces hallucinations and accurately includes the previously omitted object, 'bird', enhancing the description with essential details.
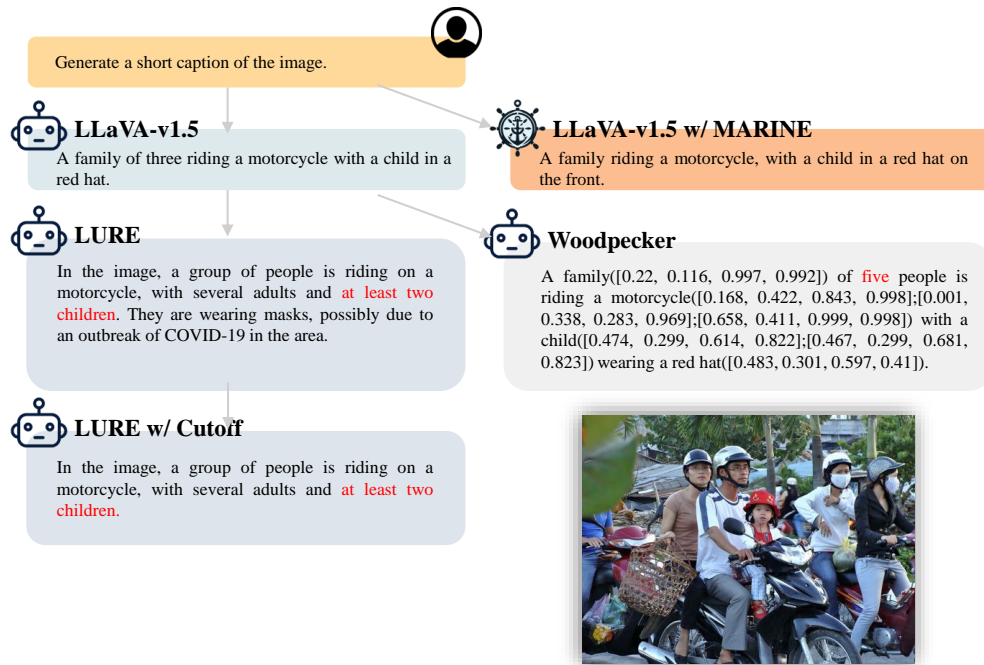
**Generate a short caption of the image.**

**LLaVA-v1.5**
A family of three riding a motorcycle with a child in a red hat.

**LLaVA-v1.5 w/ MARINE**
A family riding a motorcycle, with a child in a red hat on the front.

**LURE**
In the image, a group of people is riding on a motorcycle, with several adults and at least two children. They are wearing masks, possibly due to an outbreak of COVID-19 in the area.

**Woodpecker**
A family([0.22, 0.116, 0.997, 0.992]) of five people is riding a motorcycle([0.168, 0.422, 0.843, 0.998];[0.001, 0.338, 0.283, 0.969];[0.658, 0.411, 0.999, 0.998]) with a child([0.474, 0.299, 0.614, 0.822];[0.467, 0.299, 0.681, 0.823]) wearing a red hat([0.483, 0.301, 0.597, 0.41]).

**LURE w/ Cutoff**
In the image, a group of people is riding on a motorcycle, with several adults and at least two children.

Figure 19: A comparison of responses from baseline models and our MARINE in an image description task.



**[vqa] Based on the image, respond to this question with a short answer: Is there a toothbrush in the image?**

**MiniGPT-v2**
no

**MiniGPT-v2 w/ MARINE**
yes

**Woodpecker**
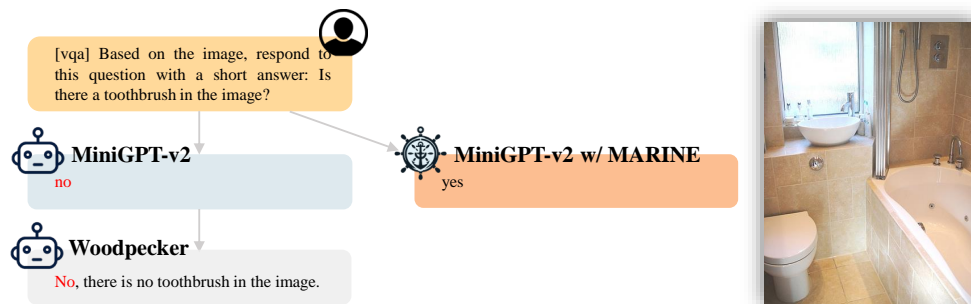No, there is no toothbrush in the image.

Figure 20: A comparison of responses from baseline models and our MARINE in POPE "yes-or-no" task. MiniGPT-v2 provides a concise response without referencing any objects. Under these circumstances, Woodpecker is unable to perform corrections via GPT-3.5 due to missing visual details. MARINE, however, successfully corrects the response while retaining MiniGPT-v2's style.

# References

ALAYRAC, J.-B., DONAHUE, J., LUC, P., MIECH, A., BARR, I., HASSON, Y., LENC, K., MENSCH, A., MILLICAN, K., REYNOLDS, M. ET AL. (2022). Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems* **35** 23716–23736.

BAZI, Y., RAHHAL, M. M. A., BASHMAL, L. and ZUAIR, M. (2023). Vision–language model for visual question answering in medical imagery. *Bioengineering* **10** 380.

BITEN, A. F., GÓMEZ, L. and KARATZAS, D. (2022). Let there be a clock on the beach: Reducing object hallucination in image captioning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*.

CARION, N., MASSA, F., SYNNAEVE, G., USUNIER, N., KIRILLOV, A. and ZAGORUYKO, S. (2020). End-to-end object detection with transformers. In *European conference on computer vision*. Springer.

CARLSSON, F., ÖHMAN, J., LIU, F., VERLINDEN, S., NIVRE, J. and SAHLGREN, M. (2022). Fine-grained controllable text generation using non-residual prompting. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

CHAMBON, P., BLUETHGEN, C., LANGLOTZ, C. P. and CHAUDHARI, A. (2022). Adapting pretrained vision-language foundational models to medical imaging domains. *arXiv preprint arXiv:2210.04133* .

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* .

CHIANG, W.-L., LI, Z., LIN, Z., SHENG, Y., WU, Z., ZHANG, H., ZHENG, L., ZHUANG, S., ZHUANG, Y., GONZALEZ, J. E., STOICA, I. and XING, E. P. (2023). Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality.

DAI, W., LI, J., LI, D., TIONG, A. M. H., ZHAO, J., WANG, W., LI, B., FUNG, P. and HOI, S. (2023a). Instructblip: Towards general-purpose vision-language models with instruction tuning.

DAI, W., LIU, Z., JI, Z., SU, D. and FUNG, P. (2023b). Plausible may not be faithful: Probing object hallucination in vision-language pre-training. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*.

DHARIWAL, P. and NICHOL, A. (2021). Diffusion models beat gans on image synthesis. *Advances in neural information processing systems* **34** 8780–8794.

DOSOVITSKIY, A., BEYER, L., KOLESNIKOV, A., WEISSENBORN, D., ZHAI, X., UNTERTHINER, T., DEHGHANI, M., MINDERER, M., HEIGOLD, G., GELLY, S. ET AL. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* .

FANG, Y., WANG, W., XIE, B., SUN, Q., WU, L., WANG, X., HUANG, T., WANG, X. and CAO, Y. (2023). Eva: Exploring the limits of masked visual representation learning at scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

FU, C., CHEN, P., SHEN, Y., QIN, Y., ZHANG, M., LIN, X., QIU, Z., LIN, W., YANG, J., ZHENG, X. ET AL. (2023). Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394* .

GAO, P., HAN, J., ZHANG, R., LIN, Z., GENG, S., ZHOU, A., ZHANG, W., LU, P., HE, C., YUE, X., LI, H. and QIAO, Y. (2023). Llama-adapter v2: Parameter-efficient visual instruction model.

GUNJAL, A., YIN, J. and BAS, E. (2023). Detecting and preventing hallucinations in large vision language models. *arXiv preprint arXiv:2308.06394* .

HO, J. and SALIMANS, T. (2021). Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications.*

HU, Z. and LI, L. E. (2021). A causal lens for controllable text generation. *Advances in Neural Information Processing Systems* **34** 24941–24955.

JI, Z., LEE, N., FRIESKE, R., YU, T., SU, D., XU, Y., ISHII, E., BANG, Y. J., MADOTTO, A. and FUNG, P. (2023). Survey of hallucination in natural language generation. *ACM Computing Surveys* **55** 1–38.

JIA, C., YANG, Y., XIA, Y., CHEN, Y.-T., PAREKH, Z., PHAM, H., LE, Q., SUNG, Y.-H., LI, Z. and DUERIG, T. (2021). Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning.* PMLR.

KAWAR, B., GANZ, R. and ELAD, M. (2022). Enhancing diffusion-based image synthesis with robust classifier guidance. *Transactions on Machine Learning Research* .

KIM, H., KIM, S. and YOON, S. (2022). Guided-tts: A diffusion model for text-to-speech via classifier guidance. In *International Conference on Machine Learning.* PMLR.

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 preprint arXiv:2311.16922* .

LI, J., LI, D., SAVARESE, S. and HOI, S. (2023a). Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597* .

LI, X., THICKSTUN, J., GULRAJANI, I., LIANG, P. S. and HASHIMOTO, T. B. (2022). Diffusion-lm improves controllable text generation. *Advances in Neural Information Processing Systems* **35** 4328–4343.

LI, X. L. and LIANG, P. (2021). Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers).*

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 preprint arXiv:2305.10355* .

LIN, S., LIU, B., LI, J. and YANG, X. (2024). Common diffusion noise schedules and sample steps are flawed. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision.*

LIN, T.-Y., MAIRE, M., BELONGIE, S., HAYS, J., PERONA, P., RAMANAN, D., DOLLÁR, P. and ZITNICK, C. L. (2014). Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*. Springer.

LIN, Z., GONG, Y., SHEN, Y., WU, T., FAN, Z., LIN, C., DUAN, N. and CHEN, W. (2023). Text generation with diffusion language models: A pre-training approach with continuous paragraph denoise. In *International Conference on Machine Learning*. PMLR.

LIN, Z., MADOTTO, A., BANG, Y. and FUNG, P. (2021). The adapter-bot: All-in-one controllable conversational model. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35.

LIU, F., LIN, K., LI, L., WANG, J., YACOOB, Y. and WANG, L. (2023a). Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565* .

LIU, F., LIN, K., LI, L., WANG, J., YACOOB, Y. and WANG, L. (2023b). Mitigating hallucination in large multi-modal models via robust instruction tuning.

LIU, H., LI, C., LI, Y. and LEE, Y. J. (2023c). Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744* .

LIU, H., LI, C., WU, Q. and LEE, Y. J. (2023d). Visual instruction tuning. In *NeurIPS*.

LOVENIA, H., DAI, W., CAHYAWIJAYA, S., JI, Z. and FUNG, P. (2023). Negative object presence evaluation (nope) to measure object hallucination in vision-language models. *arXiv preprint arXiv:2310.05338* .

LU, J., YANG, J., BATRA, D. and PARIKH, D. (2018). Neural baby talk. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

LU, P., BANSAL, H., XIA, T., LIU, J., LI, C., HAJISHIRZI, H., CHENG, H., CHANG, K.-W., GALLEY, M. and GAO, J. (2024). Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts.

MADAAN, A., SETLUR, A., PAREKH, T., POCZÓS, B., NEUBIG, G., YANG, Y., SALAKHUTDINOV, R., BLACK, A. W. and PRABHUMOYE, S. (2020). Politeness transfer: A tag and generate approach. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.

NIU, T. and BANSAL, M. (2018). Polite dialogue generation without parallel data. *Transactions of the Association for Computational Linguistics* **6** 373–389.

OUYANG, L., WU, J., JIANG, X., ALMEIDA, D., WAINWRIGHT, C., MISHKIN, P., ZHANG, C., AGARWAL, S., SLAMA, K., RAY, A. ET AL. (2022). Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* **35** 27730–27744.

PENG, N., GHAZVININEJAD, M., MAY, J. and KNIGHT, K. (2018). Towards controllable story generation. In *Proceedings of the First Workshop on Storytelling*.

PRABHUMOYE, S., BLACK, A. W. and SALAKHUTDINOV, R. (2020). Exploring controllable text generation techniques. In *Proceedings of the 28th International Conference on Computational Linguistics*.

RADFORD, A., KIM, J. W., HALLACY, C., RAMESH, A., GOH, G., AGARWAL, S., SASTRY, G., ASKELL, A., MISHKIN, P., CLARK, J. ET AL. (2021). Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR.

RIBEIRO, L. F., ZHANG, Y. and GUREVYCH, I. (2021). Structural adapters in pretrained language models for amr-to-text generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.

ROHRBACH, A., HENDRICKS, L. A., BURNS, K., DARRELL, T. and SAENKO, K. (2018). Object hallucination in image captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.

SAHARIA, C., CHAN, W., SAXENA, S., LI, L., WHANG, J., DENTON, E. L., GHASEMIPOUR, K., GONTIJO LOPES, R., KARAGOL AYAN, B., SALIMANS, T. ET AL. (2022). Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems* **35** 36479–36494.

SANCHEZ, G., FAN, H., SPANGHER, A., LEVI, E., AMMANAMANCHI, P. S. and BIDERMAN, S. (2023). Stay on topic with classifier-free guidance. *arXiv preprint arXiv:2306.17806* .

SHI, C., NI, H., LI, K., HAN, S., LIANG, M. and MIN, M. R. (2023). Exploring compositional visual generation with latent classifier guidance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

TOUVRON, H., LAVRIL, T., IZACARD, G., MARTINET, X., LACHAUX, M.-A., LACROIX, T., ROZIÈRE, B., GOYAL, N., HAMBRO, E., AZHAR, F. ET AL. (2023a). Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* .

TOUVRON, H., MARTIN, L., STONE, K., ALBERT, P., ALMAHAIRI, A., BABAEI, Y., BASHLYKOV, N., BATRA, S., BHARGAVA, P., BHOSALE, S. ET AL. (2023b). Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* .

WANG, B., WU, F., HAN, X., PENG, J., ZHONG, H., ZHANG, P., DONG, X., LI, W., LI, W., WANG, J. ET AL. (2023a). Vigc: Visual instruction generation and correction. *arXiv preprint arXiv:2308.12714* .

WANG, J., ZHOU, Y., XU, G., SHI, P., ZHAO, C., XU, H., YE, Q., YAN, M., ZHANG, J., ZHU, J. ET AL. (2023b). Evaluation and analysis of hallucination in large vision-language models. *arXiv preprint arXiv:2308.15126* .

XU, P., SHAO, W., ZHANG, K., GAO, P., LIU, S., LEI, M., MENG, F., HUANG, S., QIAO, Y. and LUO, P. (2023). Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models. *arXiv preprint arXiv:2306.09265* .

YANG, Z., WANG, J., GAN, Z., LI, L., LIN, K., WU, C., DUAN, N., LIU, Z., LIU, C., ZENG, M. ET AL. (2023). Reco: Region-controlled text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

YE, Q., XU, H., XU, G., YE, J., YAN, M., ZHOU, Y., WANG, J., HU, A., SHI, P., SHI, Y. ET AL. (2023). mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178* .

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 preprint arXiv:2310.16045* .

ZHAI, B., YANG, S., XU, C., SHEN, S., KEUTZER, K. and LI, M. (2023). Halle-switch: Controlling object hallucination in large vision language models. *arXiv e-prints* arXiv–2310.

ZHANG, H., SONG, H., LI, S., ZHOU, M. and SONG, D. (2023a). A survey of controllable text generation using transformer-based pre-trained language models. *ACM Computing Surveys* **56** 1–37.

ZHANG, Z., ZHANG, A., LI, M., ZHAO, H., KARYPIS, G. and SMOLA, A. (2023b). Multimodal chain-of-thought reasoning in language models.

ZHOU, Y., CUI, C., YOON, J., ZHANG, L., DENG, Z., FINN, C., BANSAL, M. and YAO, H. (2023). Analyzing and mitigating object hallucination in large vision-language models. *arXiv preprint arXiv:2310.00754* .

ZHU, D., CHEN, J., SHEN, X., LI, X. and ELHOSEINY, M. (2023). Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592* .