# REDUCING HALLUCINATIONS IN MULTIMODAL LARGE LANGUAGE MODELS VIA CAUSAL FUSION

## Anonymous authors

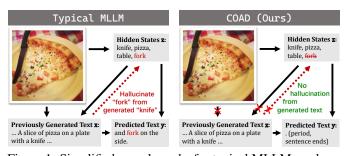
Paper under double-blind review

#### **ABSTRACT**

Multimodal Large Language Models (MLLMs) deliver detailed responses on vision-language tasks, yet remain susceptible to object hallucination (introducing objects not present in the image), undermining reliability in practice. Prior efforts often rely on heuristic penalties, post-hoc correction, or generic decoding tweaks, which do not directly intervene in the mechanisms that trigger object hallucination and thus yield limited gains. To address this challenge, we propose a causal decoding framework that applies targeted causal interventions during generation to curb spurious object mentions. By reshaping the decoding dynamics to attenuate spurious dependencies, our approach reduces false object tokens while maintaining descriptive quality. Across captioning and QA benchmarks, our framework substantially lowers object-hallucination rates and achieves state-of-the-art faithfulness without degrading overall output quality.

## 1 Introduction

Large language models (LLMs), such as GPT-4 (Achiam et al., 2023) and LLaMA (Touvron et al., 2023), have been rapidly developed and widely adopted due to their wide range of applications. To extend the capabilities of LLMs to visual tasks, multiple MLLMs have been proposed. Models such as LLaVA (Liu et al., 2024c) and MiniGPT (Zhu et al., 2023) typically project visual information into the same representational space as textual data, enabling a unified processing apthough MLLMs have shown im-



2024c) and MiniGPT (Zhu et al., Figure 1: Simplified causal graphs for typical MLLMs and our 2023) typically project visual information into the same representational space as textual data, enabling a unified processing approach via an internal LLM. Al-

pressive performance in multimodal tasks, including chatbots, visual question answering, and image captioning, they remain susceptible to *visual hallucination*.

Specifically, hallucinations in LLMs (Huang et al., 2024a) refer to cases where the model generates outputs that appear factual but are actually incorrect or ungrounded. With the introduction of visual inputs, multimodal LLMs (MLLMs) encounter a new category of hallucination: visual hallucination (Liu et al., 2024b). Visual hallucination occurs when the MLLM output diverges from the content of the input image. This undermines the reliability of the models and restricts their applicability in high-stakes real-world scenarios that demand high precision, such as medical image analysis and legal document generation.

Recently, a variety of approaches have been proposed to mitigate hallucinations in MLLMs; they can be broadly categorized into *two main strategies*: (1) The first strategy improves the model with external information, such as incorporating additional training data or retrieving knowledge from external source (Liu et al., 2024a; Yu et al., 2024; Wang et al., 2023; Chen et al., 2024a; Vu et al., 2023;

Gao et al., 2023; Varshney et al., 2023). Although these methods effectively reduce hallucinations, they often require significant effort in data collection and depend on the quality and availability of external knowledge bases. (2) The second strategy aims to reduce hallucinations without relying on additional information, instead refining the training procedures of the model or improving the attention mechanisms during inference (Yue et al., 2024; Han et al., 2024; Shi et al., 2023; Leng et al., 2023; Liu et al., 2024d; Huang et al., 2024b; Chuang et al., 2024; Deng et al., 2024; Chen et al., 2024b). However, these methods still fail to model the causal effect from visual input (e.g., images) to the generated response. They are therefore often susceptible to confounding effect or bias brought by the generated text. As a result, they tend to generate new hallucinated text based on existing hallucinated text, exacerbating hallucination.

To address these challenges, we propose Causal Object-Aware Decoding (COAD) to reduce hallucination by incorporating causal inference into the model's decoding process. Specifically, we first employ an object detector to identify visual objects in the image, delegating part of the image comprehension task to this specialized component. We then expose these structured detection results to the MLLM by finetuning the MLLM with object detection outputs as additional inputs, alongside the image and previously generated text tokens. Finally, we perform causal inference to effectively integrate the predictions from both the original pretrained model and the finetuned model to generate the response.

COAD's design improves the reliability of the MLLM via enabling targeted interventions in the model's understanding of visual objects. Furthermore, we incorporate causal inference to reduce the model's dependence on self-generated text when processing and describing images, thereby promoting more stable and less hallucinatory outputs. Our contributions are as follows:

- We formulate the generation of reliable responses as the estimation of unknown oracle predictions and introduce a new framework, dubbed Causal Object-Aware Decoding (COAD), to reduce object hallucination.
- We introduce a targeted intervention strategy that exposes and leverages visual structure, allowing the model to reason more faithfully about image content.
- We provide empirical results to demonstrate the effectiveness of our method in improving generation quality and reducing hallucination compared to state-of-the-art methods.

## 2 RELATED WORK

External Knowledge-Augmented Hallucination Mitigation. A typical strategy to mitigate hallucinations in MLLMs is to augment the model with external data. One line of work focuses on expanding or refining the training data to enhance grounding and reduce hallucinations (Liu et al., 2024a; Yu et al., 2024; Wang et al., 2023; Chen et al., 2024a). These methods typically involve curating high-quality multimodal instruction data, improving image-text alignment, or re-captioning visual content to ensure consistency with external world knowledge. By exposing the model to more reliable or better-aligned data, such approaches aim to reduce the risk of generating content that deviates from visual evidence or factual reality. Another line of research tackles hallucination at inference time by retrieving relevant information from external knowledge bases or the internet (Vu et al., 2023; Gao et al., 2023; Varshney et al., 2023). These retrieval-augmented generation methods dynamically inject grounded knowledge into the model's context, thereby improving factuality without requiring the model to memorize all details.

While both approaches have demonstrated effectiveness, they rely on either significant data curation and annotation efforts or real-time access to high-quality and up-to-date external sources. In many real-world applications, especially those involving specialized or rapidly evolving domains, such requirements may not always be feasible or reliable, highlighting the need for alternative strategies that improve factual grounding without external dependencies.

**Internal Hallucination Mitigation.** Other approaches mitigate hallucinations without relying on external data sources or retrieval mechanisms. These methods aim to improve the model's internal decision-making process by modifying its behavior during training or inference. For instance, EOS (Yue et al., 2024) encourages early stopping in sequence generation to prevent over-generation, which is often a source of factual inaccuracy. Skip-\n (Han et al., 2024) suppresses hallucinations by

skipping newline tokens, which are empirically shown to precede low-quality or fabricated continuations. Several techniques reduce the distraction caused by noisy or misleading text-conditioned inputs by selectively emphasizing attention on visual tokens. Examples include CAD (Shi et al., 2023), VCD (Leng et al., 2023), and PAI (Liu et al., 2024d), which implement visual grounding and cross-modal alignment enhancements. OPERA (Huang et al., 2024b) proposes an intervention-based decoding strategy that penalizes overconfident token predictions, which are often associated with hallucinated content. DoLa (Chuang et al., 2024) improves factual alignment by comparing generation logits from early and late transformer layers, effectively regularizing token prediction based on layer-wise consistency. In this paper, we build on this line of research by focusing on internal mechanisms to reduce hallucination, without directly relying on external knowledge bases.

## 3 METHODOLOGY

In this section, we first introduce our COAD as a causal model for the MLLM's next-token generation process, and then describe how we apply causal inference to predict the next token during inference.

## 3.1 PRELIMINARIES AND KEY INTUITION BEHIND COAD

**Problem Setting: Auto-Regressive Generation.** In this paper, we focus on auto-regressive MLLMs. At each time step, the following is given: (1) an MLLM model M; (2) an image  $\mathbf{S} \in \mathbb{R}^{c \times h \times w}$  where c is the number of channels, h is the height of the image, and w is the width of the image; (3) a sequence of previous tokens  $\mathbf{x}$  which includes both the prompt and the previously generated tokens. The model then predicts the next token y following  $\mathbf{x}$ . Specifically, the model produces a probability distribution over all possible tokens in the vocabulary, where the next token is sampled as follows:

$$y \sim P_M(y|\mathbf{x}, \mathbf{S}).$$

Here  $P_M(\cdot)$  is parameterized by the model M. In this paper, M can be a pretrained MLLM  $M_p$ , a finetuned MLLM  $M_f$ , or a hypothetical oracle MLLM  $M_*$  (more details below).

Causal Inference. Causal models provide a principled way to represent and reason about causal relationships among variables in a system (Pearl, 2009). In such models, conditional probabilities can be used to infer the distribution of hidden variables when partial observations are available. However, distinguishing causal effects from mere correlations requires the use of interventions, formalized via the do-calculus.

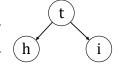


Figure 2: An exam-

ple causal graph on

temperature t, hot

Fig. 2 shows a simple causal graph where t denotes atmospheric temperature, h the sales of hot drinks, and i the sales of ice creams. The edges reflect the true causal structure: both h and i are directly influenced by t, but not by each other. In this setting, observing high hot drink sales h allows us to infer

by ice cream sales *i*.
fer sales) is also likely low.

that temperature t is likely low, which in turn implies that i (ice cream sales) is also likely low. Consequently, the conditional probability P(i|h) reflects a **misleading spurious correlation** between hot drink and ice cream sales, due to the confounder t; however, there is no causal effect from h to i.

To isolate the causal influence of h on i, we instead compute the interventional distribution P(i|do(h)), which simulates actively setting h to a fixed value while breaking its natural dependence on t. According to the rules of do-calculus, P(i|do(h)) = P(i) in this case, correctly reflecting the absence of h's causal influence on i. The formal derivation and rules can be found in (Pearl, 2009).

Causal Inference in the Context of MLLMs. A similar phenomenon arises in MLLMs' next-token prediction. Let h denote the previously generated tokens, i the predicted next token, and t the hidden states on what objects are present in the image (t is not the image itself, which could be modeled as another variable). Note that in this scenario, there should be an additional edge from h to i in the causal graph, but this distinction is irrelevant to our argument. Although h is observed during inference, t is not explicitly given, leading to **spurious correlations** in P(i|h) due to the unobserved confounder t, therefore **leading to hallucination**. To mitigate hallucination, our COAD instead considers the causal interventional distribution P(i|do(h)), which eliminates the dependence on the confounder t and predicts the next token t based only on the causal effect from t.

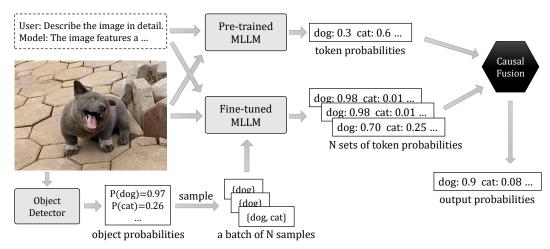


Figure 3: **Overview of our COAD.** We employ an object detector to identify the objects present in an image. The MLLM is then finetuned to condition its token predictions on both these detected objects and the input image. COAD subsequently use causal inference to combine the output distributions of both the pretrained and finetuned MLLMs to generate the final prediction.

#### 3.2 METHOD OVERVIEW

Fig. 3 illustrates the decoding (i.e., text generation) process of our COAD. It assumes access to an MLLM that can incorporate a set of detected objects as additional context during generation. To achieve this, we finetune a pretrained MLLM with object-level information. During inference, an object detector identifies likely objects in the input image and outputs a probability distribution over candidate object classes. We then sample multiple plausible object sets from this distribution.

Each sampled object set is injected as an auxiliary input into the finetuned MLLM to produce a distribution over the next token. This results in N next-token distributions, which are further combined with the distribution from the pretrained MLLM. Finally, COAD uses causal inference to combine these outputs to generate a more robust and object-aware prediction.

#### 3.3 Causal Model

Fig. 4a shows the causal model (as a causal Bayesian network) of our COAD. Given the input image S and the previous text tokens x as *observed* variables, below are key components in COAD's generative process.

**Hidden Object Variable z.** Our causal model operates at the granularity of individual token generation. At each decoding step, given the image S and the preceding (incomplete) text x, COAD infers the presence of visual objects in S through a latent binary variable  $z \in \{0,1\}^C$ , where C denotes the total number of object categories. This variable is sampled from the distribution produced by an object detector D:

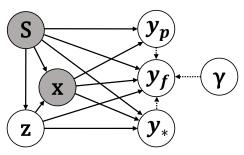
$$\mathbf{z} \sim D(\mathbf{S}),$$

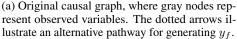
where  $D(\mathbf{S}) \in [0,1]^C$  denotes the detector's estimated probability for the presence of each object category in the image.

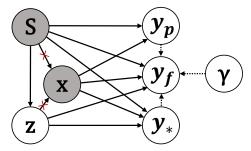
**Dual MLLMs for Generation.** To model the next-token prediction, we incorporate two MLLMs into our causal framework: a pretrained model  $M_p$  and a finetuned variant  $M_f$ . The pretrained model  $M_p$  takes as input the image S and the preceding text x, and outputs a distribution over the next token  $y_p$ . The finetuned model  $M_f$ , adapted from  $M_p$ , additionally conditions on the object variable z to produce a distribution over the next token  $y_f$ :

$$y_p \sim P_{M_p}(y_p|\mathbf{x}, \mathbf{S}),$$
  
 $y_f \sim P_{M_f}(y_f|\mathbf{x}, \mathbf{S}, \mathbf{z}),$ 

where  $P_{M_p}(y_p|\mathbf{x}, \mathbf{S})$  is the next token distribution predicted by  $M_p$ , and  $P_{M_f}(y_f|\mathbf{x}, \mathbf{S}, \mathbf{z})$  is the next token distribution predicted by  $M_f$ . In practice,  $M_p$  and  $M_f$  share most parameters for efficiency.







(b) The resulting causal graph after  $do(\mathbf{x})$  intervention. All edges going into  $\mathbf{x}$  are blocked by the intervention.

Figure 4: Illustration of our COAD's causal model before and after intervention.

**Hypothetical Oracle MLLM.** To complete the causal graph, we introduce a hypothetical oracle model  $M_*$ , which serves as an idealized reference that always produces the optimal next-token distribution. The token predicted by this oracle, denoted as  $y_*$ , is generated as follows:

$$y_* \sim P_{M_*}(y_*|\mathbf{x}, \mathbf{S}, \mathbf{z}),$$

where  $P_{M_*}(y_*|\mathbf{x}, \mathbf{S}, \mathbf{z})$  represents the oracle's ground-truth distribution conditioned on the previous text  $\mathbf{x}$ , image  $\mathbf{S}$ , and object variable  $\mathbf{z}$ .

**Mixture-Based Generation.** We hypothesize that the finetuned model  $M_f$  behaves as a mixture of the pretrained model  $M_p$  and the hypothetical oracle model  $M_*$ . At each decoding step,  $M_f$  may generate either the token predicted by  $M_p$  or the one predicted by  $M_*$ , with a certain probability. Note that this is a natural assumption:  $M_p$  serves as the initialization of  $M_f$ , and during finetuning,  $M_f$  is optimized to better approximate ground-truth signals (as represented by  $M_*$ ) while still inheriting behaviors from the original pretrained  $M_p$ .

To capture the uncertainty in  $M_f$ 's alignment between  $M_p$  and  $M_*$ , we introduce a random variable  $\gamma \in [0,1]$ , which governs the mixture proportion. This variable is drawn from a Beta distribution with hyperparameters  $\gamma_a, \gamma_b \in \mathbb{R}^+$ :

$$\gamma \sim \text{Beta}(\gamma_a, \gamma_b)$$
,

and the next-token prediction  $y_f$  is drawn approximately from a mixture of the two sources:

$$\begin{aligned} y_f &\approx \text{CategoricalMixture}(\{y_*, y_p\}, [\gamma, 1 - \gamma]) \\ &\triangleq \text{Categorical}(\gamma \times y_* + (1 - \gamma) \times y_p). \end{aligned}$$

This formulation reflects the intuition that  $M_f$  may probabilistically interpolate between following the oracle model and reverting to its pretraining prior (more details in Eqn. (2) below). It also provides an alternative generative view of  $y_f$ , which allows us to indirectly infer the oracle token  $y_*$  in Sec. 3.4.

Complete Causal Graph. Fig. 4a summarizes the causal relationships among all random variables introduced in our model. The image S and previous tokens x are the only observed variables. Note that x itself may be influenced by S and z during previous decoding steps. All other variables are conditionally generated from their respective parents according to the mechanisms described above. The dotted connections from  $y_*$ ,  $y_p$ , and  $\gamma$  to  $y_f$  indicate our hypothesis:  $M_f$  can be alternatively interpreted as a probabilistic mixture of  $M_*$  and  $M_p$ .

With the causal graph and the given observed variables, i.e., the image S and previous tokens x, our goal is to (approximately) predict the oracle next token  $y_*$  using causal inference. This will be discussed in Sec. 3.4 below.

#### 3.4 Inference Process

In this subsection, we describe how COAD employs our causal model to address the key challenges in reducing hallucinations of the MLLMs. We start by briefly discussing two key components of our

method, i.e., Causal Inference of Objects z and Estimation of Oracle Predictions, and then derive the corresponding equations that combine these two components.

Component 1: Causal Inference of Objects z. To ensure object beliefs reflect only the image content, we explicitly model them as variable z in our causal framework. Different from existing methods, where object belief is entangled in the hidden state and influenced by previous tokens  $\mathbf{x}$ , we block this dependency using an intervention  $do(\mathbf{x})$  (see Fig. 4b). This treats  $\mathbf{x}$  as externally fixed, forcing the inference of  $\mathbf{z}$  to depend solely on the image  $\mathbf{S}$  and not on prior language outputs  $\mathbf{x}$ .

Component 2: Estimation of Oracle Predictions. To approximate the oracle prediction  $y_*$ , we model the finetuned output  $y_f$  as a mixture of the pretrained model  $M_p$  and the oracle model  $M_*$ , following our assumption in Sec. 3.3. While  $y_*$  is unobservable, this mixture formulation allows us to estimate it using the available predictions  $y_f$  and  $y_p$  by  $y_f$  and  $y_p$ , respectively. This provides a principled way to bridge the gap between observed model behavior and the ideal oracle output.

Combining Components 1 & 2 to Derive the Inference Objective. By combining the previous components, our inference objective becomes computing the oracle prediction under intervention, i.e.,  $P(y_*|\mathbf{S}, \text{do}(\mathbf{x}))$ . Using Bayes' rule and standard rules of causal inference (Pearl, 2009), we have that:

$$P(y_*|\mathbf{S}, do(\mathbf{x}))$$
(1)  
=\sum\_{\mathbf{z}} P(y\_\*|\mathbf{S}, do(\mathbf{x}), \mathbf{z})P(\mathbf{z}|\mathbf{S}, do(\mathbf{x}))  
=\sum\_{\mathbf{z}} P(y\_\*|\mathbf{S}, do(\mathbf{x}), \mathbf{z})P(\mathbf{z}|\mathbf{S}) (Rule 3)  
=\sum\_{\mathbf{z}} P(y\_\*|\mathbf{S}, \mathbf{x}, \mathbf{z})P(\mathbf{z}|\mathbf{S}), (Rule 2)

This formulation rewrites the interventional query (with  $do(\cdot)$ ) using standard conditional probabilities (without  $do(\cdot)$ ), which can be estimated from observable components. We use the object detector D to compute  $P(\mathbf{z}|\mathbf{S})$ , which ensures that object beliefs are based solely on the image. The term  $P(y_*|\mathbf{S},\mathbf{x},\mathbf{z})$  represents the oracle model's prediction, which is not directly accessible. To address this, we approximate it using a mixture model. Specifically, following our hypothesized relationship between  $M_f$ ,  $M_*$ , and  $M_p$ , we have that:

$$P(y_f|\mathbf{S}, \mathbf{x}, \mathbf{z}) = \mathbb{E}_{\gamma} [\gamma P(y_*|\mathbf{S}, \mathbf{x}, \mathbf{z}) + (1 - \gamma) P(y_p|\mathbf{S}, \mathbf{x})].$$
 (2)

By rearranging Eqn. (2), we can rewrite the prediction from  $M_*$  in terms of the predictions  $y_p$  and  $y_f$  from  $M_p$  and  $M_f$ , respectively. Specifically:

$$P(y_*|\mathbf{S}, \mathbf{x}, \mathbf{z}) = \frac{1}{\mathbb{E}_{\gamma}[\gamma]} P(y_f|\mathbf{S}, \mathbf{x}, \mathbf{z}) + (1 - \frac{1}{\mathbb{E}_{\gamma}[\gamma]}) P(y_p|\mathbf{S}, \mathbf{x}) = \left(1 + \frac{\gamma_b}{\gamma_a}\right) P(y_f|\mathbf{S}, \mathbf{x}, \mathbf{z}) - \left(\frac{\gamma_b}{\gamma_a}\right) P(y_p|\mathbf{S}, \mathbf{x}).$$
(3)

**Final Inference Objective.** After substituting Eqn. (3) into Eqn. (1) and rearranging the terms, we can then rewrite our final inference objective as a combination of known quantities:

$$P(y_*|\mathbf{S}, \mathbf{do}(\mathbf{x}))$$

$$= \sum_{\mathbf{z}} P(\mathbf{z}|\mathbf{S}) [(1+\alpha) P(y_f|\mathbf{S}, \mathbf{x}, \mathbf{z}) - \alpha P(y_p|\mathbf{S}, \mathbf{x})]$$

$$= (1+\alpha) \sum_{\mathbf{z}} [P(\mathbf{z}|\mathbf{S}) P(y_f|\mathbf{S}, \mathbf{x}, \mathbf{z})] - \alpha P(y_p|\mathbf{S}, \mathbf{x}).$$
(4)

where we use the shorthand  $\alpha \triangleq \gamma_b/\gamma_a$ .

Since the dimension of  $\mathbf{z}$  can be large, directly summing over all possible object beliefs is computationally intractable. To address this, we apply Monte Carlo sampling to approximate the summation. Specifically, we draw  $\mathbf{z}_i \sim P(\mathbf{z}|\mathbf{S})$  and estimate the expectation as  $\frac{1}{N} \sum_{i=1}^n \left[ P(y_f|\mathbf{S},\mathbf{x},\mathbf{z}_i) \right]$ , where N is the number of samples. This approach offers an efficient approximation of the full marginal over  $\mathbf{z}$  while preserving the grounding of predictions in the visual input. In practice, we find that feeding the probability vector of  $\mathbf{z}$  (denoted as  $\widetilde{\mathbf{z}}$ ) directly into  $M_f$  serves as an efficient approximation of the above expectation (over  $\mathbf{z}$ ), significantly reducing computational overhead while maintaining strong performance.

Summary of COAD. To summarize, training and inference of COAD consist of the following steps:

- 1. Modify the pretrained MLLM to accept an object belief vector **z** as an additional input.
- 2. Finetune the modified MLLM using the object vectors **z** (predicted by an object detector).
- 3. At inference time, compute the next-token probability using Eqn. (4), approximating the expectation over z via Monte Carlo sampling.

Therefore, COAD enables *object-aware dehallucination* by explicitly grounding language generation in visual object beliefs. Through causal modeling and intervention, COAD ensures that predictions remain faithful to the image content, reducing reliance on spurious correlations from prior text.

# 4 EXPERIMENTS

In this section, we compare COAD with existing methods on real-world datasets.

#### 4.1 Datasets and Metrics

We use various datasets and metrics below to evaluate the MLLMs.

**POPE.** The Polling-based Object Probing Evaluation (POPE) (Li et al., 2023) employs visual question answering to assess whether an MLLM can correctly identify the presence of an object in an input image. Following the literature (Liu et al., 2024d; Huang et al., 2024b), we focus on the MSCOCO dataset with 500 images, with each image having 6 questions for each split of POPE. We evaluate the object recognition performance using the Precision, Recall, F-1, and Accuracy metrics.

**CHAIR.** Caption Hallucination Assessment with Image Relevance (CHAIR) (Rohrbach et al., 2018) is a set of widely used metrics to evaluate captioning hallucination. Following the literature (Liu et al., 2024d; Huang et al., 2024b), we use the MSCOCO dataset (Lin et al., 2014) that provides annotations for ground-truth objects in images. Specifically, CHAIR includes two metrics:

• CHAIR<sub>S</sub>, which measures the proportion of captions containing hallucinated objects relative to the total number of captions:

 $CHAIR_S = |captions|$  with hallucinated objects | / |all captions|,

• **CHAIR**<sub>I</sub>, which measures the proportion of the hallucinated objects relative to the total number of mentioned objects:

 $CHAIR_I = |hallucinated objects| / |all mentioned objects|.$ 

**MMHal-Bench.** MMHal-Bench (Sun et al., 2023) is a dataset designed to evaluate MLLMs on diverse questions where they may produce false claims about image content. These questions cover object attributes, spatial relationships, and holistic descriptions etc.

#### 4.2 Baselines

We use LLaVA-1.5-7B (Liu et al., 2024c) as the base model for all evaluated methods. For COAD, we use RTMDet (Lyu et al., 2022) as the object detector D. We compare COAD with state-of-the-art methods, including Decoding by Contrasting Layers (**DoLa**) (Chuang et al., 2024), Paying More Attention to Image (**PAI**) (Liu et al., 2024d), End-of-Sentence Decision (**EOS**) (Yue et al., 2024), Over-trust Penalty and Retrospection-Allocation (**OPERA**) (Huang et al., 2024b), Visual Contrastive Decoding (**VCD**) (Leng et al., 2023), Context-Aware Decoding (**CAD**) (Shi et al., 2023), and Object Hallucination Reduction via Adaptive Focal-Contrast Decoding (**HALC**) (Chen et al., 2024b).

#### 4.3 IMPLEMENTATION DETAILS

We finetune COAD on a subset of MSCOCO images sourced from the LLaVA dataset. To enable the model to incorporate the auxiliary input  ${\bf z}$ , we introduce a two-layer MLP projector (with hidden size 256) that maps  ${\bf z}$  into the token embedding space, following LLaVA's multimodal token integration approach. We employ LoRA ( $r=128, \alpha=256$ ), a cosine learning rate schedule with an initial learning rate of 4e-5, a batch size of 128, and train the model for 1 epoch. During inference, we use sampling by default with temperature 0.2 and a maximum of 512 output tokens. See Appendix A for more details.



**Prompt**: Generate a concise description for the image.

z-vector: (pizza: 0.89, person: 0.74, oven: 0.42, bowl: 0.36, ..., knife: 0.18, fork: 0.09, ...)

**LLaVA**: A person is cutting a pizza with a knife and fork, and the pizza is placed on a tray. The pizza has cheese and potatoes on it.

**COAD**: A pizza with one slice missing is being cut by a person wearing nail polish.

Green: Correct Objects
Red: Hallucinated Objects

Figure 5: Case study on caption generation. MSCOCO objects mentioned in the text are highlighted in red (hallucinated) or green (correct). We compare the baseline LLaVA with our COAD-enhanced model. While LLaVA hallucinates nonexistent objects (e.g., knife and fork), the z-vector produced by the object detector suggests that these objects are absent. By leveraging this signal, COAD produces a faithful caption grounded in the actual image content, consistent with the improvements shown in CHAIR metrics.

Table 1: Comparison of different methods in terms of CHAIR metrics. **Boldface** and <u>underlining</u> denote the best and the second-best performance, respectively.

Method	Base	PAI	DoLa	VCD	CAD	OPERA	EOS	HALC	COAD
$\overline{\operatorname{CHAIR}_I}\downarrow$	9.9	5.8	13.0	11.4	9.9	4.5	5.8	5.2	3.4
$\mathrm{CHAIR}_S\downarrow$	29.6	11.3	37.0	32.5	28.0	7.4	10.6	11.1	5.3

Table 2: Evaluation on MMHal-Bench across 8 hallucination dimensions: attributes (attr), adversarial objects (adv), comparison (cmp), counting (cnt), spatial relations (rel), environment (env), holistic/overall description (hol), and others (oth). **Boldface** and <u>underlining</u> denote the best and the second-best performance, respectively.

Method	Avg. Score	Hall. Rate	attr	adv	cmp	cnt	rel	env	hol	oth
Base	1.88	0.68	2.33	1.25	2.67	0.83	1.75	3.17	1.42	1.58
PAI	2.10	0.65	1.92	1.33	2.25	2.17	2.17	3.67	1.75	1.58
Dola	2.01	0.62	2.08	1.42	2.75	1.67	1.17	4.00	<u>1.75</u>	1.25
VCD	1.98	0.67	2.17	1.83	1.83	1.33	2.42	3.33	1.33	1.58
CAD	2.00	0.64	2.50	1.25	2.42	0.75	1.33	3.83	1.83	2.08
OPERA	2.09	0.65	2.58	1.67	2.67	2.50	1.58	3.08	1.17	1.50
EOS	2.08	0.62	2.67	1.33	2.67	1.00	1.83	3.17	1.58	2.42
HALC	2.12	0.64	2.33	1.67	3.00	2.25	1.67	3.42	1.33	1.33
COAD	2.52	0.52	3.58	1.83	3.33	2.08	2.08	3.50	1.33	2.42

#### 4.4 RESULTS

In this section, we compare COAD with different baselines across various datasets and metrics.

**Free-Form Generation Evaluation on CHAIR.** We first evaluate COAD on the CHAIR benchmark, which measures hallucination rates in free-form image captioning. The CHAIR benchmark includes two sub-metrics:  $CHAIR_I$  (instance-level) and  $CHAIR_S$  (sentence-level). Lower  $CHAIR_I$  and  $CHAIR_S$  indicate fewer hallucinated mentions.

As shown in Table 1, COAD achieves the best performance across all three CHAIR metrics, significantly reducing hallucinations. Specifically, it achieves 3.4 and 5.3 in terms of  $CHAIR_I$  and  $CHAIR_S$ , respectively, outperforming all existing baselines. This demonstrates that our causal object-aware decoding effectively reduces hallucination of generated captions.

Fig. 5 shows a qualitative example comparing the baseline LLaVA and our COAD. Here, LLaVA hallucinates nonexistent objects such as a *knife* and *fork*, while our COAD correctly suppresses them and generates a faithful caption. This illustrates how causal object-aware decoding helps mitigate hallucination in practice. Additional case studies are provided in Appendix E.

Table 3: POPE evaluation results on the MSCOCO dataset. **Boldface** and <u>underlining</u> denote the best and the second-best performance, respectively.

Method	Random					Popular					Adversarial				
11200100	Acc	P	R	F1	Yes	Acc	P	R	F1	Yes	Acc	P	R	F1	Yes
Base	89.0	89.3	88.6	89.0	49.6	85.0	82.6	88.7	85.5	53.7	78.8	74.0	88.8	80.8	60.0
PAI	89.3	89.6	88.9	89.2	49.6	86.1	84.2	89.0	86.5	52.9	78.9	74.4	88.3	80.7	59.4
Dola	86.3	85.3	87.7	86.5	51.4	83.0	80.5	87.1	83.6	54.1	78.2	73.9	87.4	80.1	59.2
VCD	88.8	88.8	88.7	88.8	50.0	85.4	83.6	88.1	85.8	52.7	<u>79.2</u>	74.5	88.6	81.0	59.4
CAD	88.6	88.7	88.5	88.6	49.9	84.8	82.5	88.3	85.3	53.5	78.5	74.0	87.9	80.3	59.4
OPERA	89.4	89.7	89.0	89.3	49.6	85.9	83.9	89.0	86.4	53.1	79.1	74.3	89.0	81.0	59.9
EOS	85.4	81.5	91.7	86.3	56.3	81.2	75.8	91.7	83.0	60.5	75.9	69.6	91.9	79.2	66.0
HALC	88.7	89.9	87.1	88.5	48.5	85.8	84.8	87.1	86.0	51.4	79.1	75.0	87.1	80.6	58.1
COAD	89.0	89.6	88.3	89.0	49.3	85.5	84.0	87.6	85.8	52.1	<b>79.8</b>	<b>75.8</b>	87.5	81.2	57.7

Multimodal QA Evaluation on MMHal-Bench. Table 2 shows the results on MMHal-Bench. COAD achieves the highest average score (2.52) and the lowest hallucination rate (0.52), significantly outperforming all baselines. The strong performance is consistent across multiple benchmark subsets, particularly in the Attribute, Comparison, and Relation categories, indicating improved factual accuracy and reasoning. These results further demonstrate that incorporating object-level cues effectively reduces hallucination while maintaining or enhancing generation quality.

**Object Probing Evaluation on POPE.** Table 3 shows the POPE evaluation results across three settings. COAD achieves the highest accuracy (79.8) and F1 score (81.2) on the Adversarial subset, outperforming all baselines, indicating better robustness to prompts designed to induce hallucination. In the Popular and Random subsets, it performs comparably to state-of-the-art methods in F1 while maintaining a low hallucination ratio. These results confirm that our approach effectively reduces hallucinations while preserving factual precision across diverse input types.

Ablation Studies. We conduct two ablation studies on CHAIR to better understand the source of our improvements. Specifically, we compare our full COAD with (1) "COAD ( $M_f$  Only)", which only uses the finetuned model  $M_f$  without applying our causal decoding procedure and (2) "COAD ( $\mathbf{w/o}$  z)", where we train  $M_f$  without z and perform causal decoding using this modified  $M_f$ . Table 4 shows the results. The gap between COAD and "COAD ( $M_f$  Only)" verifies the effectiveness of our causal decoding algorithm, while the gap between COAD and "COAD ( $\mathbf{w/o}$  z)" verifies the important role of z in COAD (see more discussion in Appendix B).

Table 4: Results of COAD and ablations on CHAIR. " $M_f$  only" means only using the finetuned model  $M_f$  for generation; "w/o z" means replacing  $M_f$  by a normally finetuned MLLM and applying COAD, without any z vectors involved in the whole process.

Method	$\mathrm{CHAIR}_I\downarrow\mathrm{CHAIR}_S\downarrow$							
COAD (Full)	3.4	5.3						
$COAD (M_f \text{ only})$	5.4	10.8						
COAD (w/o z)	6.9	18.1						

## 5 CONCLUSION

In this paper, we propose COAD, a novel approach to reducing hallucination in MLLMs. By combining object detection and causal inference, COAD improves the quality of generated captions and reasoning outputs. Extensive experiments on various benchmarks show that COAD consistently outperforms state-of-the-art dehallucination methods across diverse metrics and settings. Future work may include more sophisticated object representations and extend our causal modeling framework to additional multimodal tasks. Moreover, we plan to investigate the integration of temporal and spatial priors to further enhance the causal grounding of visual elements. Another promising direction is to incorporate user feedback or human-in-the-loop supervision to dynamically refine the intervention policy during inference. Finally, we aim to explore the scalability of COAD in real-world applications such as assistive vision systems and visually grounded dialogue. In terms of limitations, like many other MLLMs, maliciously manipulated inputs could affect COAD's performance. We defer a detailed discussion of limitations and potential mitigations to Appendix C.

## REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Yangyi Chen, Karan Sikka, Michael Cogswell, Heng Ji, and Ajay Divakaran. Dress: Instructing large vision-language models to align and interact with humans via natural language feedback, 2024a. URL https://arxiv.org/abs/2311.10081.
- Zhaorun Chen, Zhuokai Zhao, Hongyin Luo, Huaxiu Yao, Bo Li, and Jiawei Zhou. Halc: Object hallucination reduction via adaptive focal-contrast decoding. *arXiv preprint arXiv:2403.00425*, 2024b.
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. Dola: Decoding by contrasting layers improves factuality in large language models, 2024. URL https://arxiv.org/abs/2309.03883.
- Ailin Deng, Zhirui Chen, and Bryan Hooi. Seeing is believing: Mitigating hallucination in large vision-language models via clip-guided decoding, 2024. URL https://arxiv.org/abs/2402.15300.
- Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Y. Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu. Rarr: Researching and revising what language models say, using language models, 2023. URL https://arxiv.org/abs/2210.08726.
- Zongbo Han, Zechen Bai, Haiyang Mei, Qianli Xu, Changqing Zhang, and Mike Zheng Shou. Skip n: A simple method to reduce hallucination in large vision-language models, 2024. URL https://arxiv.org/abs/2402.01345.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, November 2024a. ISSN 1558-2868. doi: 10.1145/3703155. URL http://dx.doi.org/10.1145/3703155.
- Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. Opera: Alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation, 2024b. URL https://arxiv.org/abs/2311.17911.
- Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. Mitigating object hallucinations in large vision-language models through visual contrastive decoding, 2023. URL https://arxiv.org/abs/2311.16922.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13*, pp. 740–755. Springer, 2014.
- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Mitigating hallucination in large multi-modal models via robust instruction tuning, 2024a. URL https://arxiv.org/abs/2306.14565.
- Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li, and Wei Peng. A survey on hallucination in large vision-language models, 2024b. URL https://arxiv.org/abs/2402.00253.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2024c. URL https://arxiv.org/abs/2310.03744.

- Shi Liu, Kecheng Zheng, and Wei Chen. Paying more attention to image: A training-free method for alleviating hallucination in lvlms, 2024d. URL https://arxiv.org/abs/2407.21771.
- Chengqi Lyu, Wenwei Zhang, Haian Huang, Yue Zhou, Yudong Wang, Yanyi Liu, Shilong Zhang, and Kai Chen. Rtmdet: An empirical study of designing real-time object detectors, 2022. URL https://arxiv.org/abs/2212.07784.
- Judea Pearl. Causality. Cambridge university press, 2009.
- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. *arXiv preprint arXiv:1809.02156*, 2018.
- Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Scott Wen tau Yih. Trusting your evidence: Hallucinate less with context-aware decoding, 2023. URL https://arxiv.org/abs/2305.14739.
- Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. URL https://arxiv.org/abs/2302.13971.
- Neeraj Varshney, Wenlin Yao, Hongming Zhang, Jianshu Chen, and Dong Yu. A stitch in time saves nine: Detecting and mitigating hallucinations of Ilms by validating low-confidence generation, 2023. URL https://arxiv.org/abs/2307.03987.
- Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and Thang Luong. Freshllms: Refreshing large language models with search engine augmentation, 2023. URL https://arxiv.org/abs/2310.03214.
- Lei Wang, Jiabang He, Shenshen Li, Ning Liu, and Ee-Peng Lim. Mitigating fine-grained hallucination by fine-tuning large vision-language models with caption rewrites, 2023. URL https://arxiv.org/abs/2312.01701.
- Qifan Yu, Juncheng Li, Longhui Wei, Liang Pang, Wentao Ye, Bosheng Qin, Siliang Tang, Qi Tian, and Yueting Zhuang. Hallucidoctor: Mitigating hallucinatory toxicity in visual instruction data, 2024. URL https://arxiv.org/abs/2311.13614.
- Zihao Yue, Liang Zhang, and Qin Jin. Less is more: Mitigating multimodal hallucination from an eos decision perspective, 2024. URL https://arxiv.org/abs/2402.14545.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models, 2023. URL https://arxiv.org/abs/2304.10592.

## A IMPLEMENTATION DETAILS

We finetune COAD on a subset of MSCOCO images sourced from the LLaVA dataset. To enable the model to incorporate the auxiliary input  $\mathbf{z}$ , we introduce a two-layer MLP projector (with hidden size 256) that maps  $\mathbf{z}$  into the token embedding space, following LLaVA's multimodal token integration approach. We employ LoRA ( $r=128, \alpha=256$ ), a cosine learning rate schedule with an initial learning rate of 4e-5, a batch size of 128, and train the model for 1 epoch. To mitigate the model's dependence on prior context, we apply Gaussian noise ( $\sigma=0.005$ ) to the embeddings of previous tokens with a probability of 0.5 during training. During inference, we use sampling by default with temperature 0.2 and a maximum of 512 output tokens. In implementing Eqn. (4), we find that it is more effective to perform fusion in the logit space rather than in the probability space. Therefore, we replace  $P(y_f|\mathbf{S},\mathbf{x},\mathbf{z})$  and  $P(y_p|\mathbf{S},\mathbf{x})$  with their corresponding logits before computing the fused output, which is subsequently converted back to the probability space via softmax.

All experiments were conducted on a single machine with 8 NVIDIA RTX A5000 GPUs (24GB each), an AMD EPYC 7282 16-Core Processor (64 threads), and 256GB RAM. Finetuning typically took around 16 hours per model. Caption generation on 5,000 images took between 30 minutes and 2 hours, depending on the generation length.

For evaluation, we use sampling to generate outputs for all baseline methods, except for OPERA. Since OPERA is built on top of beam search, its outputs are generated using beam search with a beam size of 3 instead. For the hyperparameter  $\alpha$  in COAD, we set it to 1.5 for text generation tasks (CHAIR and MMHal-Bench) and 0.1 for POPE.

#### B FURTHER ANALYSIS OF ABLATION STUDIES

Effect of Finetuning and Effectiveness of Our Causal Decoding Algorithm. Since COAD involves finetuning an MLLM, a natural question is whether the observed gains are simply due to finetuning rather than our proposed causal decoding algorithm. To examine this, we directly evaluate the finetuned model  $M_f$  without applying our decoding strategy. As shown in Table 1,  $M_f$  alone achieves only part of the improvements, indicating that finetuning by itself cannot account for the performance of COAD and verifying the effectiveness of our causal decoding algorithm.

**Role of the Vector z.** Another question is whether the improvements come merely from contrasting  $M_f$  and  $M_p$  during causal fusion, regardless of our vector  $\mathbf{z}$ . To verify this, we remove  $\mathbf{z}$  when training  $M_f$  (i.e., a standard finetuning setting) and then apply our causal decoding procedure using this variant of  $M_f$ . The results from Table 1 show a clear drop compared to COAD, demonstrating that  $\mathbf{z}$  plays an essential role in enabling effective causal fusion.

### C LIMITATIONS AND FUTURE WORK

**Dependence on Finetuning.** Our approach currently requires a LoRA finetuning step for adaptation. While this is practical in many settings, it reduces the plug-and-play convenience of COAD. We experimented with a training-free variant that injects the causal vector directly as a prompt, which already yields strong improvements on POPE benchmarks (e.g., F1-Rand = 95.4, F1-Pop = 90.0, F1-Adv = 85.8), outperforming all baselines. However, this variant is more sensitive to detector errors and less robust on captioning tasks. These results nonetheless highlight the generality of COAD and suggest promising directions for reducing computational cost, such as improving detector reliability or designing dedicated training methods that allow a single MLLM to simulate the causal signal.

**Domain Mismatch.** COAD relies on detectors trained on specific distributions. If the test image domain diverges significantly, the causal signal may become insufficient. One direction is to investigate zero-shot or domain-adaptive detectors to mitigate this issue.

**Adversarial Vulnerability.** Maliciously manipulated inputs or detector outputs could affect COAD's performance. However, its modular design allows for safeguard components (e.g., adversarial detection at the detector level), which we leave as future extensions.

**Residual Text Priors.** In rare cases with extremely strong linguistic priors, causal interventions may not fully suppress hallucinations. In rare cases with extremely strong linguistic priors, causal interventions may not fully suppress hallucinations. Future improvements may involve designing stronger intervention mechanisms or complementary signals that better counteract such priors.

**Scope and Generality.** Our study mainly targets object hallucinations. Broader validation on other types of hallucinations and across more MLLMs is a promising next step, facilitated by COAD's general token-based interface.

#### D USE OF LARGE LANGUAGE MODELS

We used a large language model (ChatGPT, OpenAI, 2025) for correcting typos and minor grammatical issues. No parts of the technical content, experiments, or analysis were generated by the model.

## E MORE QUALITATIVE EXAMPLES

 We provide more qualitative examples in Fig. 6 and Fig. 7, where MSCOCO objects mentioned in the text are highlighted in red (hallucinated) or green (correct).



**z-vector:** [person: 0.93, handbag: 0.76, cell phone: 0.72, potted plant: 0.5, clock: 0.26, bottle: 0.24, **bench: 0.2**, ...]

**LLaVA:** A woman is sitting on a bench, holding a cell phone and talking to another woman. They are both wearing sandals.

**COAD:** A woman in a pink shirt holds a cell phone, while another woman in a dress sits next to her.



**z-vector:** [toilet: 0.86, cat: 0.84, book: 0.46, **dog: 0.17**, ...]

**LLaVA:** A dog is standing on a toilet, sniffing the toilet bowl. The toilet lid is open, and the dog appears to be curious about the water inside.

**COAD:** A cat is sitting on the toilet and appears to be drinking water from the toilet bowl.



**z-vector:** [train: 0.91, chair: 0.26, bird: 0.21, **bench: 0.2**, person: 0.16, clock: 0.15, traffic light: 0.14, ...]

**LLaVA:** A yellow and white train is traveling down the tracks, pulling into a station. The train is positioned next to a building, and there is a bench nearby for passenger to sit. The scene is captured on November 1, 2012.

**COAD:** A train is stopped at a train station, and there is a yellow line on the platform.

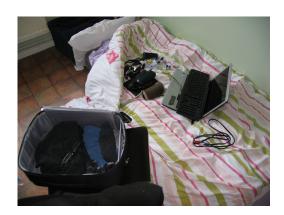


**z-vector:** [banana: 0.87, apple: 0.87, orange: 0.83, bowl: 0.33, **dining table: 0.28**, wine glass: 0.11, cup: 0.05, ...]

**LLaVA:** A clear glass plate is filled with a variety of fruits, including apple, orange, and a banana. The plate is placed on a table, and the fruits are arranged in a visually appealing manner.

**COAD:** A clear glass plate holds a variety of fruits, including banana, apple, pears, and an orange.

Figure 6: Case studies on hallucination reduction.



**z-vector:** [laptop: 0.88, bed: 0.86, suitcase: 0.63, keyboard: 0.44, mouse: 0.42, backpack: 0.36, handbag: 0.33, **cell phone: 0.31**, remote: 0.15, person: 0.14, ...]

**LLaVA:** A bed with a laptop and a cell phone on it. The bed is unmade and has a striped blanket. There is also a suitcase on the bed, and a pair of glasses nearby.

**COAD:** A laptop is sitting on top of a bed next to a suitcase.



**z-vector:** [bicycle: 0.89, person: 0.83, motorcycle: 0.75, backpack: 0.34, bottle: 0.27, handbag: 0.18, ..., **chair: 0.05**, ...]

**LLaVA:** A man is working on a motorcycle in a garage, standing on a red stool to reach the motorcycle. He is wearing jeans and a black shirt. The motorcycle is parked on a lift, and there are other motorcycle in the background.

**COAD:** A man is working on a motorcycle, lifting the front wheel off the ground using a lift.



**z-vector:** [hot dog: 0.91, person: 0.77, sandwich: 0.31, bottle: 0.16, pizza: 0.07, ..., **dining table: 0.04**, ...]

**LLaVA:** A person is holding a hot dog with cheese and ketchup on a bun. The hot dog is placed on a table, and the person's hand is visible in the foreground.

**COAD:** A hot dog with ketchup and cheese being held by a person.



**z-vector:** [person: 0.85, bird: 0.83, cup: 0.83, sandwich: 0.64, **bench: 0.55**, dining table: 0.5, fork: 0.33, chair: 0.32, knife: 0.22, ...]

**LLaVA:** A bird is standing on a plate with a half-eaten sandwich, which is placed on a dining table. The bird seems to be interested in the sandwich, possibly trying to get a bite. The scene takes place near a body of water, with a bench nearby.

**COAD:** A half-eaten sandwich sits on a plate with ketchup, and a bird is standing nearby, possibly interested in the remaining food.

Figure 7: Case studies on hallucination reduction.