# Causal Decoding for Hallucination-Resistant Multimodal Large Language Models

**Shiwei Tan[1*], Hengyi Wang[1*], Weiyi Qin[1*], Qi Xu[2*], Zhigang Hua[2*], Hao Wang[1]**

**[*]Equal Contribution, [1]Rutgers University, [2]Meta Ranking AI**

## 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 approach via an internal LLM. Although MLLMs have shown impressive performance in multimodal tasks, including chatbots, visual question answering, and image captioning, they remain susceptible to *visual hallucination*.
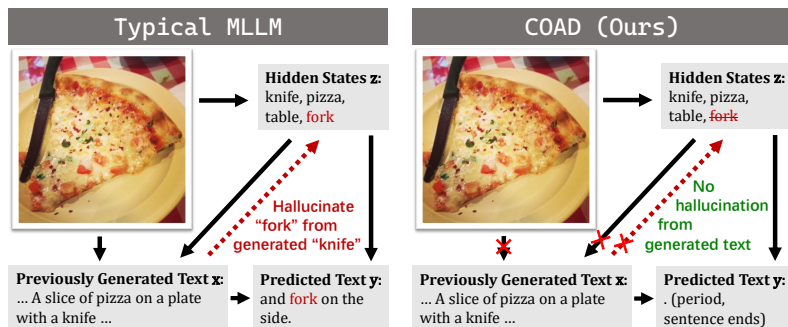


Figure 1: Simplified causal graphs for typical MLLMs and our COAD. **Left:** Typical MLLMs implicitly hallucinate objects (e.g., "fork") in the hidden states $\mathbf{z}$ due to previously generated text $\mathbf{x}$ (e.g., "knife"). **Right:** Our COAD performs causal inference to remove links between the hidden states $\mathbf{z}$ and generated text $\mathbf{x}$, thereby avoiding 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., 2024a; 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 Yan & Wang (2023); Wang et al. (2020); Pearl (2009) 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; this is inspired by hierarchical Bayesian deep learning Wang & Yeung (2016; 2020); Wang et al. (2024) and the causality literature Pearl (2009); Wang et al. (2020); Yan & Wang (2023). 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 object 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. CLIP-guided decoding (Deng et al., 2024b) reduces hallucination by incorporating a CLIP-based image-text consistency score into a sentence-level beam search, adjusting the beam scores beyond the MLLM's own likelihood. 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.** We consider an auto-regressive MLLM that, at each decoding step, receives: (1) a model $M$; (2) an input 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; and (3) a sequence of previous tokens $\mathbf{x}$ (including the prompt and generated tokens so far). The model predicts the next token $y$ by sampling from

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

Here $M$ may be a pretrained MLLM $M_p$, a finetuned model $M_f$, or a hypothetical oracle $M_*$ (introduced later).

**Causal Inference.** Causal models provide a principled framework for distinguishing true causal effects from correlations induced by confounders (Pearl, 2009). A central tool is the use of interventions, denoted by $\mathrm{do}(\cdot)$, which remove spurious dependencies when analyzing the effect of one variable on another.
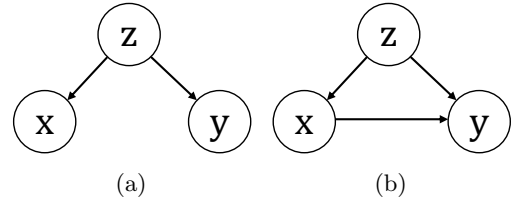


Figure 2: **Illustration of confounding.** (a) $z$ induces a spurious association between $x$ and $y$ even without a causal effect. (b) Adding $x \rightarrow y$ introduces a genuine causal effect, but $P(y|x)$ remains confounded by $z$.

**Illustrative Example.** We next illustrate the effect of a confounder using the two causal graphs in Figure 2. Let $z$ denote temperature, $x$ the hot drink sales, and $y$ the ice cream sales.

In Figure 2a, the edges reflect the causal structure where both $x$ and $y$ are influenced by the confounder $z$, but there is no direct causal effect from $x$ to $y$. In this setting, observing high hot drink sales $x$ implies that temperature $z$ is likely low, which in turn suggests that $y$ (ice cream sales) is also likely low. Consequently, even though $x$ does not causally affect $y$, the conditional probability $P(y|x)$ differs from $P(y)$. This discrepancy reflects a *misleading spurious correlation* introduced by the confounder $z$.

To isolate the causal influence of $x$ on $y$, we instead compute the interventional distribution $P(y|\mathrm{do}(x))$, which simulates actively setting $x$ to a fixed value while breaking its natural dependence on $z$. According to the rules of do-calculus, $P(y|\mathrm{do}(x)) = P(y)$ in this case, correctly reflecting the absence of $x$'s causal influence on $y$. The formal derivation and rules can be found in (Pearl, 2009).

Figure 2b then adds a direct causal edge $x \to y$ on top of Figure 2a. In this case, hot drink sales $x$ indeed have a direct causal impact on ice cream sales $y$ (e.g., some customers may buy a cold item after consuming a hot drink). However, the conditional probability $P(y|x)$ still *overestimates* this causal effect, because $z$ continues to influence both $x$ and $y$. The confounder $z$ thus causes $P(y|x)$ to capture both the genuine causal effect of $x$ on $y$ and the additional dependence mediated by $z$, motivating the use of interventional quantities such as $P(y|\mathrm{do}(x))$ to isolate the true causal influence.

**Analogy to MLLM Next-Token Prediction.** The structure in Figure 2b directly corresponds to MLLM decoding. Let $\mathbf{x}$ denote the previously generated tokens, $y$ the next token to be predicted, and $\mathbf{z}$ the model's hidden states representing its belief about which objects are present in the image (note that $\mathbf{z}$ is not the image itself, which could be modeled separately). Because $\mathbf{z}$ is a backdoor variable connecting to both $\mathbf{x}$ and $y$, the conditional probability $P(y|\mathbf{x})$ still **overestimates** the causal effect from $\mathbf{x}$ to $y$, i.e., overestimates the likelihood of certain tokens $y$, therefore potentially leading to hallucination.

**Key Intuition of COAD.** To mitigate object hallucination, COAD explicitly models $\mathbf{z}$ as a variable representing object beliefs and replaces the standard conditional distribution with the interventional one $P(y|\mathrm{do}(\mathbf{x}), \mathbf{z})$. This removes the spurious dependence introduced by the confounder $\mathbf{z}$ and predicts the next token $y$ based solely on the true causal effects of $\mathbf{x}$ and $\mathbf{z}$.

**Generative vs. Recognition Causal Models.** The causal perspective adopted in this work follows the *recognition/inference* view commonly used in recent causal analyses of vision-language models (Mao et al., 2021; 2022). Instead of modeling how the physical world generates an image (which would typically imply a generative direction such as $\mathbf{z} \to \mathbf{S}$), our objective is to describe how an MLLM processes an observed image $\mathbf{S}$ and pre-existing text $\mathbf{x}$ in order to form internal beliefs and predict the next token.

Under this recognition view, $\mathbf{z}$ denotes the model's internal belief about which objects are present in the input image, not a latent variable causing the image itself. The information flow from $\mathbf{S}$ to $\mathbf{z}$ therefore represents the model's inference procedure, consistent with modern architectures such as LLaVA, where visual features are computed by the vision encoder before any interaction with textual tokens (see Appendix C). This viewpoint provides the conceptual grounding for the causal structures introduced in the following sections.

### 3.2 Formal Definition of Object Hallucination

Let $\mathbf{S}$ be the input image, $\mathbf{z}^*$ the ground-truth set of visual objects, and $p_\theta(y|\mathbf{x}, \mathbf{S})$ the model's predictive distribution for the next token $y$ given previous tokens $\mathbf{x}$ and the image $\mathbf{S}$. Let $p^*(y|\mathbf{x}, \mathbf{z}^*)$ denote the "ground-truth" conditional distribution, i.e., the distribution produced by an ideal model that fully respects the true visual semantics of the image.

We define object hallucination as the divergence between these two distributions:

$$D(p_\theta(y|\mathbf{x}, \mathbf{S}) \,\|\, p^*(y|\mathbf{x}, \mathbf{z}^*)),$$

where $D(\cdot\|\cdot)$ may be a KL or another divergence measure. A large divergence indicates that the model assigns high probability to tokens that contradict the true visual object content, thereby producing object hallucination.

In existing multimodal LLMs, hallucination often arises because the hidden states $\mathbf{z}$ may encode nonexistent objects based on previous tokens $\mathbf{x}$ rather than on the true image content $\mathbf{S}$. This can lead the model to generate tokens $y$ that are not visually grounded. We empirically validate this phenomenon for LLaVA using a linear-probe analysis of its internal object-existence beliefs; see Appendix B for details.

To address this issue, our COAD (i) uses a detector-derived proxy $\hat{\mathbf{z}}$ to approximate the visual constraints $\mathbf{z}^*$ in the ideal distribution $p^*(y|\mathbf{x}, \mathbf{z}^*)$, and (ii) intervenes on both the internal object-related hidden states
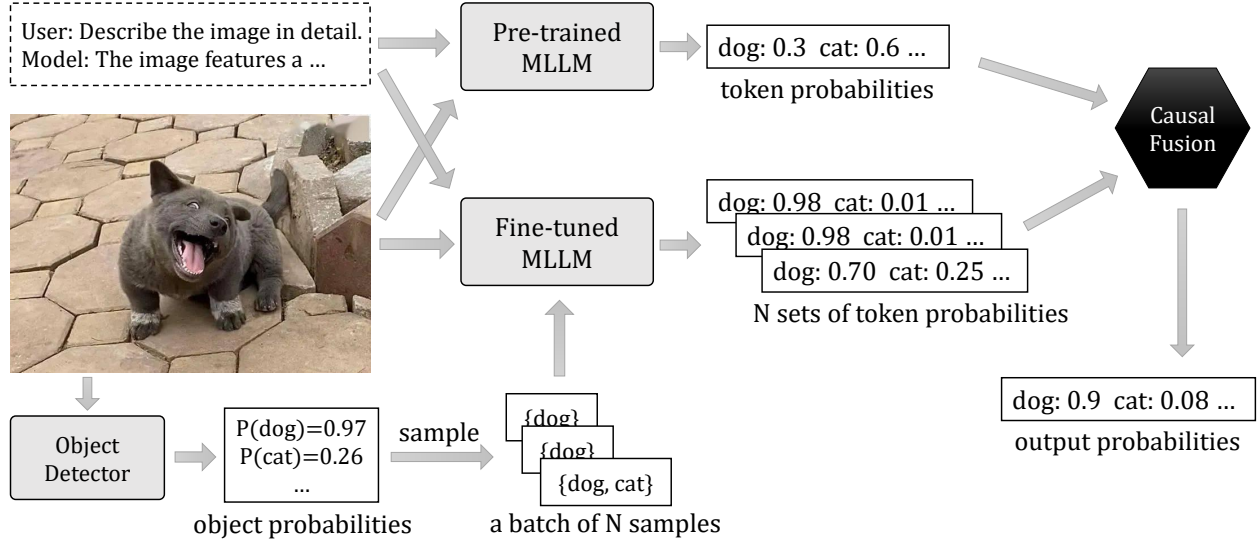
Figure 3: **Overview of our COAD.** We employ an object detector to identify the objects present in an image. The MLLM is then finetuned to condition its token predictions on both these detected objects and the input image. COAD subsequently use causal inference to combine the output distributions of both the pretrained and finetuned MLLMs to generate the final prediction.

$\mathbf{z}$ and the previous tokens $\mathbf{x}$ to block the spurious information pathway from $\mathbf{x}$ to $y$. These interventions reduce the divergence above and therefore mitigate object hallucination.

### 3.3 Method Overview

Figure 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.

### 3.4 Causal Model of a Standard MLLM

Before introducing the causal model underlying COAD, we first describe the temporal causal structure of a standard MLLM during autoregressive decoding. This view makes explicit how information flows across timesteps and clarifies the relationship between the image, the evolving text sequence, and the model's internal object-related variables.

**Temporal Structure.** The input image $\mathbf{S}$ and the initial prompt $\mathbf{x}^{(0)}$ remain fixed throughout decoding. The variable $\mathbf{z}$ represents the model's belief about which objects are present in the image and is determined solely by $\mathbf{S}$; it therefore remains constant across timesteps. Here, $\mathbf{z}$ is a conceptual variable capturing image-conditioned object beliefs, rather than a token-level hidden state that evolves during decoding (see Appendix A for further discussion). This treatment is consistent with the architecture of modern vision-language models such as LLaVA, where visual features are produced by the vision encoder before any interaction with textual tokens (see Appendix C). The only time-varying variables are the evolving text sequence $\mathbf{x}^{(t)}$ and the next token $y^{(t)}$.

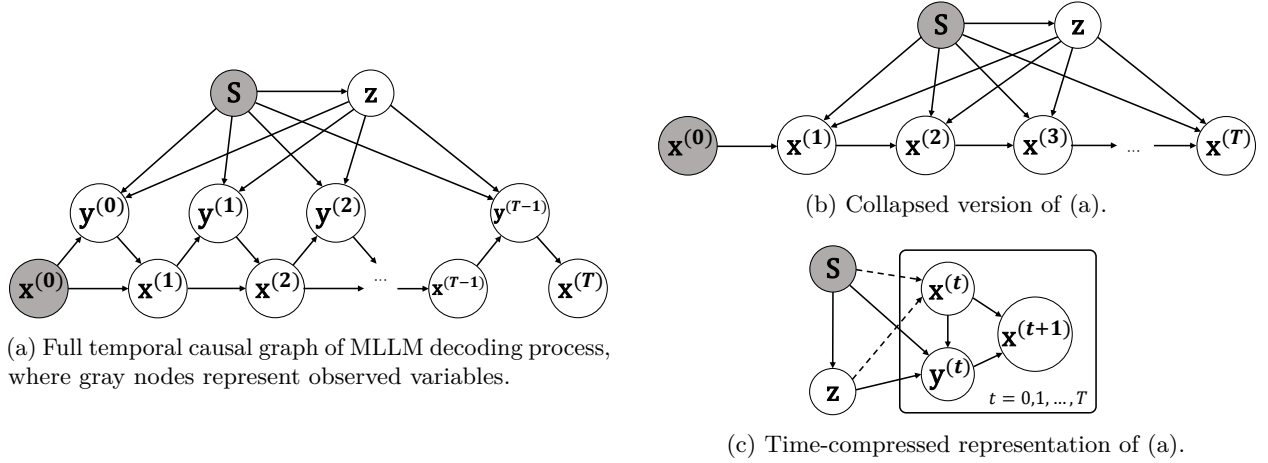Figure 4a illustrates the resulting temporal causal graph, which contains the following edges:

(a) Full temporal causal graph of MLLM decoding process, where gray nodes represent observed variables.

(b) Collapsed version of (a).

(c) Time-compressed representation of (a).

Figure 4: **Rolled-out causal structures of the MLLM decoding process. (a)** Full temporal (rolled-out) causal graph over decoding timesteps: the image $\mathbf{S}$ and initial prompt $\mathbf{x}^{(0)}$ are fixed, $\mathbf{z}$ denotes the image-derived object-belief variable (constant over time), and $\mathbf{x}^{(t)}$, $y^{(t)}$ evolve according to $\mathbf{x}^{(t)}, \mathbf{S}, \mathbf{z} \to y^{(t)}$ and $\mathbf{x}^{(t)}, y^{(t)} \to \mathbf{x}^{(t+1)}$. **(b)** Collapsed version of (a) obtained by treating $y^{(t)}$ as an implicit step in the autoregressive update, so that each $\mathbf{x}^{(t)}$ ($t > 0$) is jointly determined by $\mathbf{S}$, $\mathbf{z}$, and $\mathbf{x}^{(t-1)}$. **(c)** Time-compressed representation in which the accumulated influence of all previous timesteps on $\mathbf{x}^{(t)}$ is summarized by dashed edges from $(\mathbf{S}, \mathbf{z})$ to $\mathbf{x}^{(t)}$.

- $\mathbf{S} \to \mathbf{z}$, indicating that the object-belief variable $\mathbf{z}$ is inferred from the image.

- $\mathbf{x}^{(t)}, \mathbf{S}, \mathbf{z} \to y^{(t)}$, reflecting that each predicted token depends on the current text, the visual input, and the inferred object beliefs.

- $\mathbf{x}^{(t)}, y^{(t)} \to \mathbf{x}^{(t+1)}$, reflecting the autoregressive update where each new token $y^{(t)}$ is appended to the old sequence $\mathbf{x}^{(t)}$ to form the new input sequence $\mathbf{x}^{(t+1)}$, i.e., the new input sequence is $\mathbf{x}^{(t+1)} = [\mathbf{x}^{(t)}, y^{(t)}]$.

**Collapsed Representation.** If the intermediate variables $y^{(t)}$ are treated as implicit steps in the autoregressive update, the subgraph of $(\mathbf{x}^{(t)}, y^{(t)}, \mathbf{x}^{(t+1)})$ can be collapsed into a single transition $\mathbf{x}^{(t)} \to \mathbf{x}^{(t+1)}$, producing the structure shown in Figure 4b. In this form, each $\mathbf{x}^{(t)}$ ($t > 0$) is jointly determined by $\mathbf{S}$, $\mathbf{z}$, and $\mathbf{x}^{(t-1)}$.

**Time-Compressed View.** Compressing the time axis of Figure 4a yields the representation in Figure 4c, which summarizes the cumulative influence of past timesteps on $\mathbf{x}^{(t)}$ via dashed edges from $(\mathbf{S}, \mathbf{z})$ to $\mathbf{x}^{(t)}$. This compressed form closely matches the causal structure underlying a single decoding step, and provides the foundation upon which we build the causal model for COAD in the next subsection.

### 3.5 Causal Model of COAD

Figure 5a shows the causal model (as a causal Bayesian network) of our COAD. Given the input image $\mathbf{S}$ and the previous text tokens $\mathbf{x}$ as *observed* variables, below are key components in COAD's generative process.

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

**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 global prior Beta distribution with hyperparameters $\gamma_a, \gamma_b \in \mathbb{R}^+$, which are fixed across dataset:

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

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

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

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 Equation 2 below). It also provides an alternative generative view of $y_f$, which allows us to indirectly infer the oracle token $y_*$ in Section 3.6.

**Complete Causal Graph.** Figure 5a summarizes the causal relationships among all random variables introduced in our model. The image $\mathbf{S}$ and previous tokens $\mathbf{x}$ are the only observed variables. Note that $\mathbf{x}$ itself may be influenced by $\mathbf{S}$ and $\mathbf{z}$ 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 $\mathbf{S}$ and previous tokens $\mathbf{x}$, our goal is to (approximately) predict the oracle next token $y_*$ using causal inference. This will be discussed in Section 3.6 below.
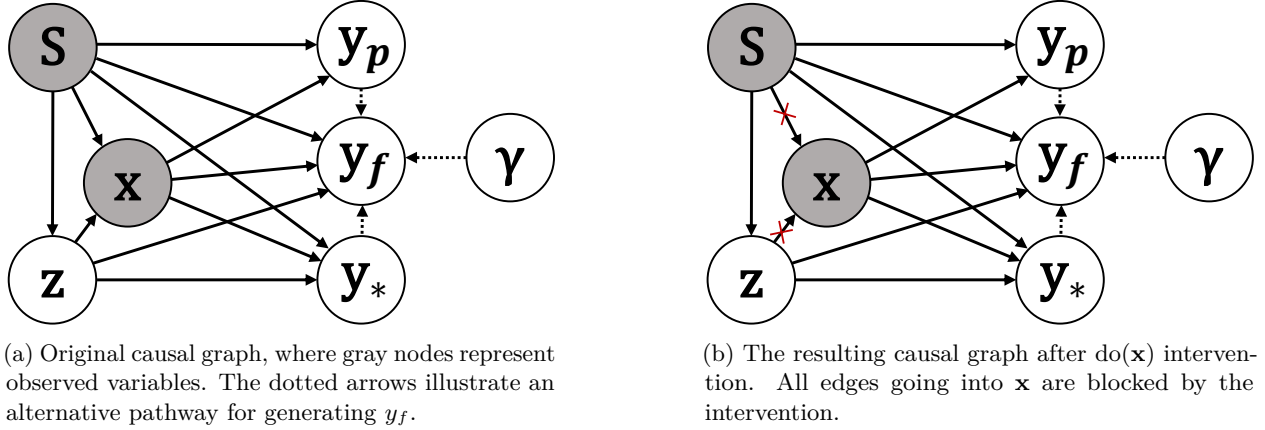
(a) Original causal graph, where gray nodes represent observed variables. The dotted arrows illustrate an alternative pathway for generating $y_f$.

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

Figure 5: Illustration of our COAD's causal model before and after intervention.

### 3.6 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 $\mathbf{z}$ 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 Figure 5b). 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 Section 3.5. While $y_*$ is unobservable, this mixture formulation allows us to estimate it using the available predictions $y_f$ and $y_p$ by $M_f$ and $M_p$, 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:

$$P(y_*|\mathbf{S}, \text{do}(\mathbf{x})) \tag{1}$$
$$= \sum_{\mathbf{z}} P(y_*|\mathbf{S}, \text{do}(\mathbf{x}), \mathbf{z}) P(\mathbf{z}|\mathbf{S}, \text{do}(\mathbf{x}))$$
$$= \sum_{\mathbf{z}} P(y_*|\mathbf{S}, \text{do}(\mathbf{x}), \mathbf{z}) P(\mathbf{z}|\mathbf{S}) \qquad \text{(Rule 3)}$$
$$= \sum_{\mathbf{z}} P(y_*|\mathbf{S}, \mathbf{x}, \mathbf{z}) P(\mathbf{z}|\mathbf{S}), \qquad \text{(Rule 2)}$$

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:

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

By rearranging Equation 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:

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

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

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

where we use the shorthand $\alpha \triangleq \gamma_b/\gamma_a$. Since only the ratio $\alpha = \gamma_b/\gamma_a$ appears in the final expression, the Beta distribution's parameters $(\gamma_a, \gamma_b)$ (as hyperparameters) effectively have only one degree of freedom during inference. In practice, We therefore treat $\alpha$ as a single global hyperparameter (see Appendix A for additional details). Notably, computing the closed-form expression in Equation 4 is equivalent to sampling many values of $\gamma \sim \mathrm{Beta}(\gamma_a, \gamma_b)$ and averaging their contributions in expectation. This closed-form solution makes COAD both more efficient and conceptually consistent with the underlying mixture interpretation.

Since the dimension of $\mathbf{z}$ can be large, directly summing over all possible object-belief vectors is computationally intractable. We consider two practical implementations to approximate the expectation over $\mathbf{z}$ in Equation 4: (1) *Monte Carlo sampling*, where we treat the detector's output $\widetilde{z} \in [0, 1]^C$ as the parameters of $C$ Bernoulli distributions and sample $N$ binary vectors $z_i \in \{0, 1\}^C$ (where $i = 1, 2, \ldots, N$) from them, so that the term $\sum_{\mathbf{z}} P(\mathbf{z}|\mathbf{S})P(y_f|\mathbf{S}, \mathbf{x}, \mathbf{z})$ is approximated by $\frac{1}{N}\sum_{i=1}^{N} P(y_f|\mathbf{S}, \mathbf{x}, \mathbf{z}_i)$; and (2) *probability-based approximation*, where we directly feed the probability vector $\widetilde{\mathbf{z}}$ into $M_f$, which empirically provides an efficient approximation to the same expectation while avoiding sampling. Due to its high efficiency, we adopt the second implementation (although it is an approximation) for all our experiments.

**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 $\mathbf{z}$ as an additional input.

2. Finetune the modified MLLM using the object vectors $\mathbf{z}$ (predicted by an object detector).

3. At inference time, compute the next-token probability using Equation 4, approximating the expectation over $\mathbf{z}$ 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$_S$**, which measures the proportion of captions containing hallucinated objects relative to the total number of captions:

$$\text{CHAIR}_S = |\text{captions with hallucinated objects}| \ / \ |\text{all captions}|,$$

- **CHAIR$_I$**, which measures the proportion of the hallucinated objects relative to the total number of mentioned objects:

$$\text{CHAIR}_I = |\text{hallucinated objects}| \ / \ |\text{all mentioned objects}|.$$

**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. The benchmark spans eight hallucination dimensions that assess different aspects of visual grounding: (1) *Object attribute*: the ability to correctly describe properties such as color or shape; (2) *Adversarial object*: the ability to recognize when a queried object is not present rather than hallucinating it; (3) *Comparison*: the ability to compare attributes or properties across multiple objects; (4) *Counting*: the ability to estimate the number of referred objects; (5) *Spatial relation*: the ability to reason about relative positions among objects; (6) *Environment*: the ability to infer aspects of the surrounding scene or background; (7) *Holistic description*: the ability to provide accurate global descriptions of the entire image; (8) *Others*: the ability to recognize text or symbols and reason based on observable visual information. Scores for each dimension are computed as the proportion of responses deemed non-hallucinatory. Higher scores across these dimensions indicate better grounding and fewer hallucinations. Following the MMHal-Bench evaluation protocol (Sun et al., 2023), we use GPT-4 to perform this judgment.

## 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 $D$.

As discussed in Section 1, existing hallucination-mitigation approaches fall into two broad categories: (1) external-knowledge-based methods and (2) internal architecture/decoding modifications. Since COAD belongs to the second family and does not rely on external data retrieval, all baselines evaluated here are drawn from this internal-mechanism category to ensure a fair comparison. Under this setting, we compare COAD with state-of-the-art internal-mechanism 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 **z**, we introduce a two-layer MLP projector (with hidden size 256) that maps **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 6: **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 underlining denote the best and the second-best performance, respectively.

| Method | Base | PAI | DoLa | VCD | CAD | OPERA | EOS | HALC | COAD |
|--------|------|-----|------|-----|-----|-------|-----|------|------|
| $\text{CHAIR}_I \downarrow$ | 9.9 | 5.8 | 13.0 | 11.4 | 9.9 | <u>4.5</u> | 5.8 | 5.2 | **3.4** |
| $\text{CHAIR}_S \downarrow$ | 29.6 | 11.3 | 37.0 | 32.5 | 28.0 | <u>7.4</u> | 10.6 | 11.1 | **5.3** |

## 4.4 Main 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: $\text{CHAIR}_I$ (instance-level) and $\text{CHAIR}_S$ (sentence-level). Lower $\text{CHAIR}_I$ and $\text{CHAIR}_S$ 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 $\text{CHAIR}_I$ and $\text{CHAIR}_S$, respectively, outperforming all existing baselines. This demonstrates that our causal object-aware decoding effectively reduces hallucination of generated captions.

Figure 6 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 more faithful caption.[1] This illustrates how causal object-aware decoding helps mitigate hallucination in practice. Additional case studies are provided in Appendix F.

**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

---

[1]While COAD successfully removes the hallucinated objects such as the "knife" and "fork", its output also includes the phrase "one slice missing", which may be an inaccurate description of the pizza. This type of error concerns the *attributes* of an already-present object rather than the presence of additional objects. Since COAD is designed specifically to mitigate *object hallucination*, such attribute-level inconsistencies fall outside the scope of what our method targets.

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 underlining 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** | 1.75 | 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** |

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

| Method | Random | | | | | Popular | | | | | Adversarial | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 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 | 79.2 | 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 | **79.8** | **75.8** | 87.5 | **81.2** | 57.7 |

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 (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 (w/o **z**)" verifies the important role of **z** in COAD (see more discussion in Appendix D).

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 | CHAIR$_I \downarrow$ | CHAIR$_S \downarrow$ |
|---|---|---|
| COAD (Full) | **3.4** | **5.3** |
| COAD ($M_f$ only) | 5.4 | 10.8 |
| COAD (w/o **z**) | 6.9 | 18.1 |

### 4.5 Runtime and Computational Overhead

Beyond hallucination metrics, we also compare the computational overhead of COAD with existing methods. To ensure consistency with the CHAIR evaluation, we select 100 images from the MSCOCO subset used in our CHAIR experiments and ask each method to generate free-form descriptions of the images using the

Table 5: Decoding throughput of different methods on 100 MSCOCO images used in the CHAIR evaluation. We report the number of tokens per second during generation (higher is better).

| Method | Base | COAD (Ours) | OPERA | PAI | DoLa | VCD | CAD | EOS | HALC |
|---|---|---|---|---|---|---|---|---|---|
| #tokens/s ↑ | 24.37 | 10.49 | 4.52 | 43.62 | 29.38 | 7.98 | 7.92 | 9.91 | 7.32 |

same decoding setting as in Section 4.4. All methods are evaluated on a single GPU. We measure decoding throughput in tokens per second, which reflects the effective per-token computational cost of each method. Table 5 shows the results.

**Detector Overhead.** COAD invokes the object detector only once per image before decoding. Since autoregressive decoding dominates the total computational cost for MLLMs, this one-time detection overhead is relatively small: the RTMDet detector used in our implementation processes each image in about 0.10 seconds, which is negligible compared with the cumulative cost of token generation on long outputs.

**Dual-Model Decoding.** COAD evaluates both the pretrained model and the object-aware finetuned model at each decoding step. When executed sequentially on a single GPU, this dual-model decoding results in roughly half the throughput of the base LLaVA model (10.49 vs. 24.37 tokens/s in Table 5). However, the two forward passes are independent and can be executed fully in parallel on different GPUs, allowing throughput close to that of single-model decoding in practical multi-GPU deployments.

Compared to other hallucination-mitigation methods, COAD remains computationally competitive. In particular, it is significantly faster than multi-step refinement and beam-search/backtracking approaches such as OPERA (4.52 tokens/s), and comparable to other decoding-modification methods like VCD, CAD, EOS, and HALC. Overall, COAD achieves strong reductions in hallucination with a moderate and well-characterized runtime overhead.

## 5    Conclusion

In this paper, we propose COAD, a novel approach to reducing object 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. Since the number of object categories is determined by the detector, it would also be interesting to explore open-vocabulary detectors (e.g., GLIP (Li et al., 2022)), which may allow COAD to operate over a substantially richer and more flexible object space. 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. Another limitation is that COAD primarily targets object hallucination; extending the causal modeling framework to other forms of hallucination, such as attribute, relational, and global-scene inconsistencies, remains an important direction for future work. We defer a detailed discussion of limitations and potential mitigations to Appendix E.

## 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, 2024a. URL https://arxiv.org/abs/2402.15300.

Ailin Deng, Zhirui Chen, and Bryan Hooi. Seeing is believing: Mitigating hallucination in large vision-language models via clip-guided decoding. *arXiv preprint arXiv:2402.15300*, 2024b.

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.

Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10965–10975, 2022.

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.

Chengzhi Mao, Augustine Cha, Amogh Gupta, Hao Wang, Junfeng Yang, and Carl Vondrick. Generative interventions for causal learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3947–3956, 2021.

Chengzhi Mao, Kevin Xia, James Wang, Hao Wang, Junfeng Yang, Elias Bareinboim, and Carl Vondrick. Causal transportability for visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7521–7531, 2022.

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 llms 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.

Hao Wang and Dit-Yan Yeung. Towards bayesian deep learning: A framework and some existing methods. *TDKE*, 28(12):3395–3408, 2016.

Hao Wang and Dit-Yan Yeung. A survey on bayesian deep learning. *CSUR*, 53(5):1–37, 2020.

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.

Yibin Wang, Haizhou Shi, Ligong Han, Dimitris Metaxas, and Hao Wang. Blob: Bayesian low-rank adaptation by backpropagation for large language models. In *NeurIPS*, 2024.

Yuhao Wang, Vlado Menkovski, Hao Wang, Xin Du, and Mykola Pechenizkiy. Causal discovery from incomplete data: A deep learning approach. In *AAAI StarAI Workshop*, 2020.

Jingquan Yan and Hao Wang. Self-interpretable time series prediction with counterfactual explanations. In *ICML*, 2023.

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  Details on Implementation and Causal Graphs

**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 Equation 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.

**Causal Graphs and Hidden States.** There are two different hidden states:

- The **image-based, static hidden state $\mathbf{z}$**, which corresponds to the output of the MLLM's vision encoder (i.e., $H_v$ that we will further explain in Appendix C). In a causal MLLM such as LLaVA, the input follows the structure **"image | user prompt | generated tokens"**, and the hidden state of a token position is affected only by tokens that appear **before** it. Therefore, here $\mathbf{z}$ is determined solely by the image and aligns most closely with our interpretation of the object belief variable $\mathbf{z}$ in our causal graphs in Figure 4 and Figure 5. *It is not influenced by $\mathbf{x}$.*

- The **response-based, dynamic hidden state $\mathbf{h}$ (or $\mathbf{h}^{(t)}$)**, which corresponds to the hidden states when the MLLM generates the $t$-th response token. These hidden states are influenced by $\mathbf{x}$.

We further clarify the role of the dynamic hidden state $\mathbf{h}^{(t)}$ in the causal graph shown in Figure 7. In Figure 7a:

- the image $\mathbf{S}$ generates the image-based, static hidden state $\mathbf{z}$,

- $\mathbf{S}$, $\mathbf{z}$, and $\mathbf{x}^{(t)}$ then jointly generate (influence) the response-based, dynamic hidden state $\mathbf{h}^{(t)}$, and

- $\mathbf{h}^{(t)}$ then influences the generated token $y^{(t)}$.

As shown in Figure 7b, we can actually collapse $h^{(t)}$ and the *red* part of Figure 7a to have an equivalent causal graph in Figure 7b, which matches our causal graphs in Figure 4 and Figure 5.

(a) Casual graph with hidden states $\mathbf{h}^{(t)}$. The hidden states that decide $y^{(t)}$ are affected by $\mathbf{S}$, $\mathbf{z}$, and $\mathbf{x}^{(t)}$.

(b) Collapsing hidden states $\mathbf{h}^{(t)}$ into the generation of $y^{(t)}$. The resulting graph corresponds to the causal graph of COAD in Figure 5.
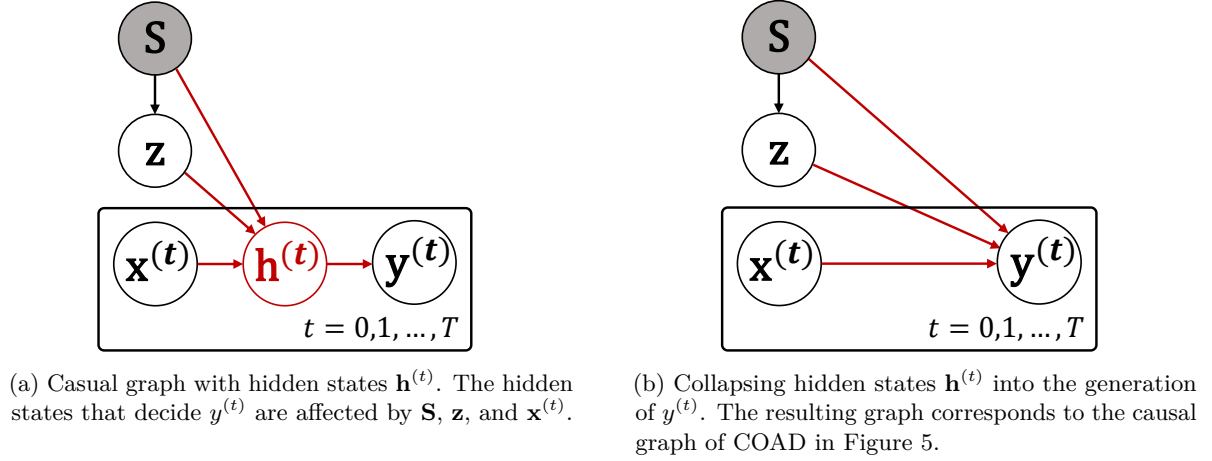
Figure 7: **The role of hidden states in the decoding causal graph of an MLLM.** This figure illustrates a causal graph that explicitly includes token-level hidden states $\mathbf{h}^{(t)}$ during generation **(a)**, and shows that collapsing these hidden states into the generation of $y^{(t)}$ yields an equivalent causal graph consistent with the one used by COAD **(b)**.

## B Evidence of Hallucinatory Object Beliefs in LLaVA

To examine whether a multimodal LLM (MLLM) can form *incorrect internal beliefs* about object existence (corresponding to an incorrect $\mathbf{z}$ in our formulation), we conduct a linear-probe experiment on LLaVA. The goal is to determine whether LLaVA forms internal object-existence beliefs that deviate from the actual content of the image.

### B.1 Experimental Setup

We train a linear classifier to probe whether LLaVA internally believes each object category to be present in the image. The classifier

- takes as **input** the embedding for each token, and

- outputs a $C$-dimensional vector $\mathbf{o} \in [0,1]^C$ indicating the estimated probability that each of the $C$ objects exists in the image.

Here, $\mathbf{o}$ reflects the MLLM's object-existence belief in its hidden states, which we mentioned in Appendix A.

**Input to the Linear Probe.** For each token in the dialog sequence, LLaVA produces a hidden state after every transformer decoder layer. We concatenate these hidden states together with the embedding before the first layer to obtain the probe input. Thus, each token in the prompt-response sequence contributes one training sample for the linear probe.

**Conversation Sampling.** We collect dialogs generated by LLaVA using MSCOCO images and the prompt "Please describe every object in this image in detail." After LLaVA completes the response, we extract hidden states for all tokens (including image tokens, query tokens, and generated response tokens) to construct probe input samples.

**Target Labels.** For each image, we run LLaVA once with the same prompt and use a rule-based method to determine whether each object category is mentioned in the resulting text description. These detected mentions serve as labels for training the linear probe. The same label is assigned to all probe samples within a given image-prompt-response tuple.

**Discussion on the Relation between Input to the Linear Probe and the Hidden States.** In a causal MLLM such as LLaVA, the input follows the structure **"image | user prompt | generated tokens"**, and the hidden state of a token position is affected only by tokens that appear **before** it. Therefore:

- when the input to the linear probe is the hidden state in the "image" part, the output **o** serves as an estimation of the **image-based, static hidden state z**;

- when the input to the linear probe is the hidden state in the "generated tokens" part, the output **o** serves as an estimation of the **response-based, dynamic hidden state h (or $\mathbf{h}^{(t)}$)**.

### B.2 Findings

We apply the trained linear probe to a sampled conversation and visualize the predicted object probabilities over time (Figure 8).

- The input image contains a person riding a skateboard with **no benches** in the background.

- The probe output in Figure 8(a) indicates that LLaVA internally assigns a high probability (i.e., brighter, yellow color in the heatmap) to a nonexistent *bench* around certain image-token positions.

- Figure 8(b) shows LLaVA's probed probability of the object *"bench"* over the image tokens (patches), demonstrating high *"bench"* probabilities in some image tokens (highlighted in yellow); this indicates that the baseline LLaVA's hallucinated belief may originate from the image-perception stage.

- Figure 8(c) shows that LLaVA hallucinates the object *"bench"* in the response.

- In contrast, COAD produces a response that does not mention any benches, as shown in Figure 8(d).

### B.3 Conclusion

This experiment demonstrates that LLaVA can form incorrect internal beliefs about object existence, i.e., it may estimate an incorrect **z** for certain objects, which subsequently leads to hallucinations. These results provide direct empirical evidence for the assumption illustrated in Figure 1 of the main paper.

## C LLaVA Architecture and the Formation of Visual Features

To clarify why the object-belief variable **z** is treated as time-invariant in the causal model of a standard MLLM (Section 3.4), we briefly describe the relevant components of the LLaVA architecture. A schematic diagram is shown in Figure 9.
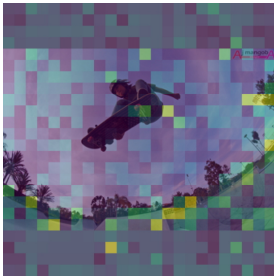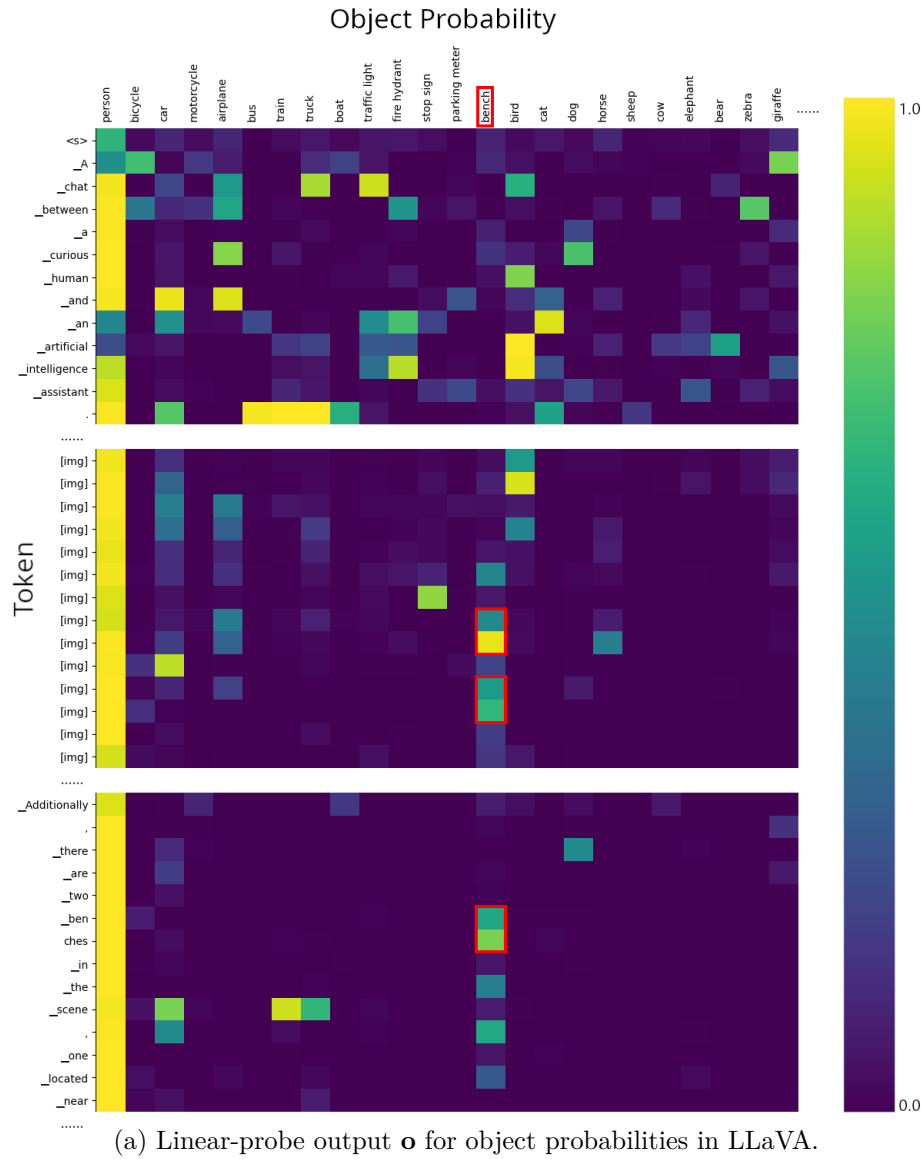
### C.1 Visual Feature Extraction Before Language Interaction

In LLaVA, the input image is first processed by a vision encoder (e.g., CLIP's ViT backbone) to produce a set of image features, denoted $Z_v$. These features are then passed through a projector layer to obtain the final visual embeddings $H_v$, which are fed into the language model.

Importantly:

- Both the vision encoder and the projector operate *independently* of the language model.

- No textual tokens are involved when computing $Z_v$ or $H_v$.

- The visual embeddings $H_v$ remain fixed throughout decoding.

The projected embeddings $H_v$ are inserted into the token sequence as "image tokens," and only after this insertion does the multimodal Transformer start attending jointly over image tokens and text tokens.

(a) Linear-probe output **o** for object probabilities in LLaVA.



**LLaVA**: The image captures a skateboarder in mid-air, performing a trick on a ramp. The skateboarder is wearing a helmet for safety while executing the jump. The skateboard is ... Additionally, there are two benches in the scene, one located near the left edge of the image and the other towards the right side.

**COAD**: In the image, a young person is performing a skateboard trick, flying through the air while riding the skateboard. The skateboard is clearly visible beneath the skateboarder, and the person appears to be in mid-air, showcasing their skill. The scene captures the excitement and athleticism of the sport.

(b) Visualization of probed probability of the object *"bench"* over all image tokens.

(c) LLaVA's response hallucinates the object *"benches"*.

(d) COAD eliminates the bench-related hallucination.

Figure 8: **Visualization of LLaVA's hallucinatory object beliefs.** This figure provides empirical evidence that LLaVA can form incorrect internal object-existence beliefs (corresponding to an incorrect **z**) and shows how COAD corrects the resulting hallucination. **(a)** Linear-probe output of LLaVA's object probabilities for all tokens over the full dialog, showing that the probed probability for the object *"bench"* increases both around *certain image tokens* and around the *tokens where LLaVA actually generates the word "benches"* (in red boxes). **(b)** The baseline LLaVA's probed probability of the object *"bench"* over the image tokens (patches), showing high *"bench"* probabilities in some image tokens (highlighted in yellow); this indicates that the hallucinated belief may originate from the image-perception stage. **(c)** LLaVA hallucinates the object *"bench"* in the response. **(d)** COAD's response, which eliminates the bench hallucination.
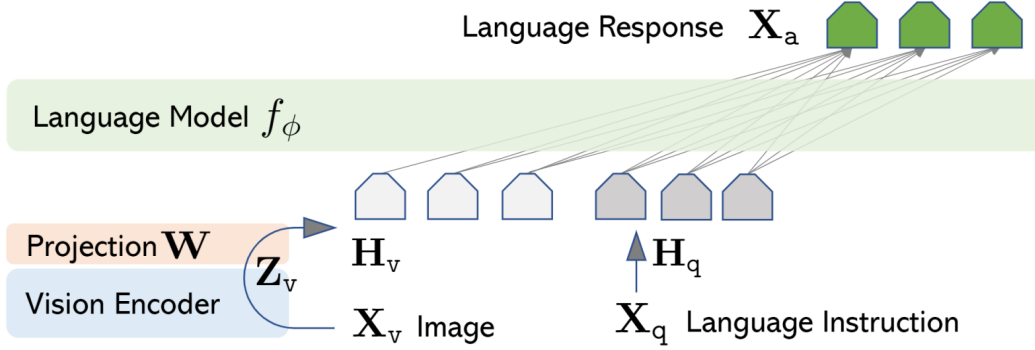
Figure 9: **Simplified LLaVA architecture (Liu et al., 2024c).** The image is encoded by a vision encoder and then projected into the language model's embedding space. The resulting visual embeddings $H_v$ are inserted as image tokens. Because $H_v$ is computed independently of textual tokens, the object-belief variable **z** corresponding to $H_v$ remains fixed during decoding. (This figure is adapted from (Liu et al., 2024c).)

### C.2  Relation to the Object-Belief Variable z

In our causal formulation, the variable **z** represents the model's internal belief about which objects exist in the image. Conceptually, **z** corresponds to the portion of the visual embeddings $H_v$ that encode object presence or absence. Since $H_v$ is generated solely from the image (via the vision encoder and projector) and does not depend on previously generated text tokens, **z** is naturally treated as a **static** variable. This matches the assumption used in Section 3.4, where only $\mathbf{x}^{(t)}$ and $y^{(t)}$ evolve over time.

## D  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 **z**. To verify this, we remove **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 **z** plays an essential role in enabling effective causal fusion.

## E  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.

## F   More Qualitative Examples

We provide more qualitative examples in Figure 10 and Figure 11, 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 10: 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:** [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:** [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:** [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 11: Case studies on hallucination reduction.