# PAS : Prelim Attention Score for Detecting Object Hallucinations in Large Vision–Language Models

Nhat Hoang-Xuan[1,2]    Minh Vu[1]    My T. Thai[2]    Manish Bhattarai[1]

[1]Los Alamos National Laboratory, Los Alamos, NM, USA    [2]University of Florida, Gainesville, FL, USA

xhoang@lanl.gov    mvu@lanl.gov    mythai@cise.ufl.edu    ceodspspectrum@lanl.gov

## Abstract

*Large vision–language models (LVLMs) are powerful, yet they remain unreliable due to object hallucinations. In this work, we show that in many hallucinatory predictions the LVLM effectively ignores the image and instead relies on previously generated output ("prelim") tokens to infer new objects. We quantify this behavior via the mutual information between the image and the predicted object conditioned on the prelim, demonstrating that weak image dependence strongly correlates with hallucination. Building on this finding, we introduce the Prelim Attention Score (PAS), a lightweight, training-free signal computed from attention weights over prelim tokens. PAS requires no additional forward passes and can be computed on the fly during inference. Exploiting this previously overlooked signal, PAS achieves state-of-the-art object-hallucination detection across multiple models and datasets, enabling real-time filtering and intervention.*

## 1. Introduction

Following the success of Large Language Models (LLMs), Large Vision-Language Models (LVLMs) brought visual perception and natural language reasoning together. They achieve new frontiers on tasks including but not limited to spatial reasoning [8], medical visual question answering [35], and visual document understanding [21]. However, hallucination remains a persistent problem, and multimodal capabilities come with a new type of hallucination: *object hallucination*, in which the model mentions objects that are not present in the image [28]. Thus, detecting object hallucination in LVLMs is crucial for ensuring their credibility and reliability, especially in practical scenarios requiring factual correctness [29].

Many existing approaches to object hallucination detection rely on annotations/references [22, 28] or external judge models [18, 34, 36]. These approaches are limited by the lack of references in many real-world cases and

the potential unreliability of judge models [27]. Training-free, reference-free approaches avoid this problem by relying only on *intrinsic* information (inputs/outputs and model internals) for detection.

Effectively utilizing the limited intrinsic information from the LVLM inference process is vital for training-free, reference-free approaches. These methods have considered various sources, such as output logits [26, 37] or image token attention [17]. Notably, they tend to focus on most informative tokens, namely image tokens, while crucially overlooking the low-information preliminary output tokens (which we call the prelim). In contrast, we argue that the model's dependence on this low-information part of the input can signify unreliability, and the inclusion of this information can help with object hallucination detection. We develop our hypothesis with a novel information-theoretic formulation and support our hypothesis with evidence (Fig. 3, Tab. 1). From our findings, we propose PAS to quantify the dependence of an object token on the prelim, which effectively acts as an object hallucination detector. PAS is based on attention weights computed during model inference, hence its overhead is minimal. Figure 1 illustrates our findings and our proposed method. Our primary contributions are presented below:

- We investigate and promote the role of *preliminary output tokens* in determining object hallucination, a source of information overlooked by previous methods when constructing a detector.
- We theorize that excessive focus on prelim tokens can be linked to an unreliable, more hallucinatory alternative operating mode of the LVLM. We provide an information-theoretic formalization to validate our hypothesis.
- Based on the above observations, we propose PAS, an attention-based score that achieves state-of-the-art object hallucination detection performance across models and datasets.

**Organization.** In Section 2, we review existing object hallucination detection methods for LVLM and related lines of research. Section 3 elaborates on LVLM token genera-

1

tion dependencies and formalizes the detection problem. In Section 4, we theorize the role of the prelim in object hallucination and develop an information-theoretic method to quantify this for detection. We also propose PAS score, an alternative efficient attention-based approach. Section 5 details the realization of both approaches. Section 6 presents experiments, ablations, and analyses. Finally, Section 7 concludes and discusses limitations.

## 2. Related Work

**Object Hallucination Detection.** Object hallucination is a phenomenon where the model outputs refer to objects that are inconsistent with or do not exist in the input images [22, 27, 28]. Existing approaches to object hallucination detection can be classified based on the information that they use:

- **Logits-based approaches** [26, 37] primarily rely on just the outputs of the LVLM for hallucination detection. They generally compute some kind of information-theoretic uncertainty measure based on probabilistic interpretation of the token logits. While principled, since they use the least information (compared to other categories), they tend to have worse detection performance.
- **Representation-based approaches** [16, 27] typically focus on image tokens and their compatibility with predicted objects. They leverage Logit Lens [1] to map intermediate image hidden states to output probabilities and use this as a proxy for image-object compatibility.
- **Attention-based approaches** [17] leverage the attention weights to quantify the dependency of output tokens on input tokens. However, similar to representation-based approaches, they mostly focus on image tokens.

Existing approaches tend to focus on image tokens and overlook information in the prelim tokens. In contrast, we observe that frequently when hallucination occurs, the model exhibits high dependency on prelim tokens. This observation is detailed in Section 4.1. Based on this observation, we develop a detector that takes the prelim into account, and obtain an efficient hallucination detection that achieves state-of-the-art performance.

**Related research directions.** Apart from *internal* information, some object hallucination methods leverage *external* information, such as an external oracle [10, 18, 34, 36]. In many cases, this reliance can be problematic since bigger models are costly to run and the external models themselves can be wrong. In this work, we focus on the setting where no external information is used and no additional training is performed, and compare against baselines with the same assumptions. Furthermore, some existing works focus on *sentence*-level or *segment*-level hallucination in LVLMs [13, 20, 34], while object hallucination is *token*-level, and thus those approaches cannot be directly applied

to this task. Finally, a line of research [3, 15, 19, 25] focuses on *mitigation* of object hallucination, which directly alters the model, its inputs, or the generation process in order to reduce hallucination. Mitigation methods do not directly lead to detection and also do not eliminate all hallucinations, hence our work and mitigation methods are complementary. We provide additional discussion of related works in the Appendix.

## 3. Preliminary

**Notation.** Let $\Phi$ be the LVLM. Let $\mathbf{v}, \mathbf{t}$ denote the image tokens and instruction tokens, respectively. Let $\mathbf{x} = (\mathbf{v}, \mathbf{t})$ be the input tokens, and $\mathbf{y}$ be the output tokens. Note that we omit the BOS token in $\mathbf{x}$ for brevity. Let $m = |\mathbf{x}|$ be the number of input tokens, and $n = |\mathbf{x}| + |\mathbf{y}|$ be the total number of tokens.

**Token generation dependency.** In this work, we focus on LVLMs which are autoregressive decoders. For autoregressive LVLMs, the probability of producing the output sequence $\mathbf{y}$ is given as:

$$\Pr(\mathbf{y} \mid \mathbf{x}) = \prod_{k=m+1}^{n} \Pr(y_k \mid \mathbf{y}_{<k}, \mathbf{x}) \qquad (1)$$

where

$$\Pr(\cdot \mid \mathbf{y}_{<k}, \mathbf{x}) = \mathrm{softmax}\left(\Phi(\cdot \mid \mathbf{y}_{<k}, \mathbf{x})\right) \qquad (2)$$

represents the output of a single forward pass of the model $\Phi$. From Eq. (1), the generation of $y_k$ depends on the input $\mathbf{x} = (\mathbf{v}, \mathbf{t})$ and *previously* generated tokens $\mathbf{y}_{<k}$ (which we call the *prelim*), therefore investigating these components can help us understand hallucinatory predictions.

As discussed in Section 2, existing approaches tend to extract signals from the image tokens $\mathbf{v}$ since it contains the "factual" information that the model should focus on, while the prelim contains incomplete information [15]. Conversely, in this work, we show that strong signals for object hallucination detection can be extracted from the prelim.

**Attention mechanism.** The attention mechanism can be used to quantify the dependency of the generation of a token on previous tokens. This is because many popular LVLMs use Large Language Model (LLM) backbones that have decoder-only Transformer architecture (e.g, LLaMA [32], QwenLM [5]). In this architecture, an LLM consists of only self-attention and MLP layers. Since the MLP layers operate independently on each token, the attention layers are the only path for information exchange between tokens in a single forward pass. Therefore, looking at attention weights can give insight on the information flow. In particular, looking at the attention weights of an output token can reveal
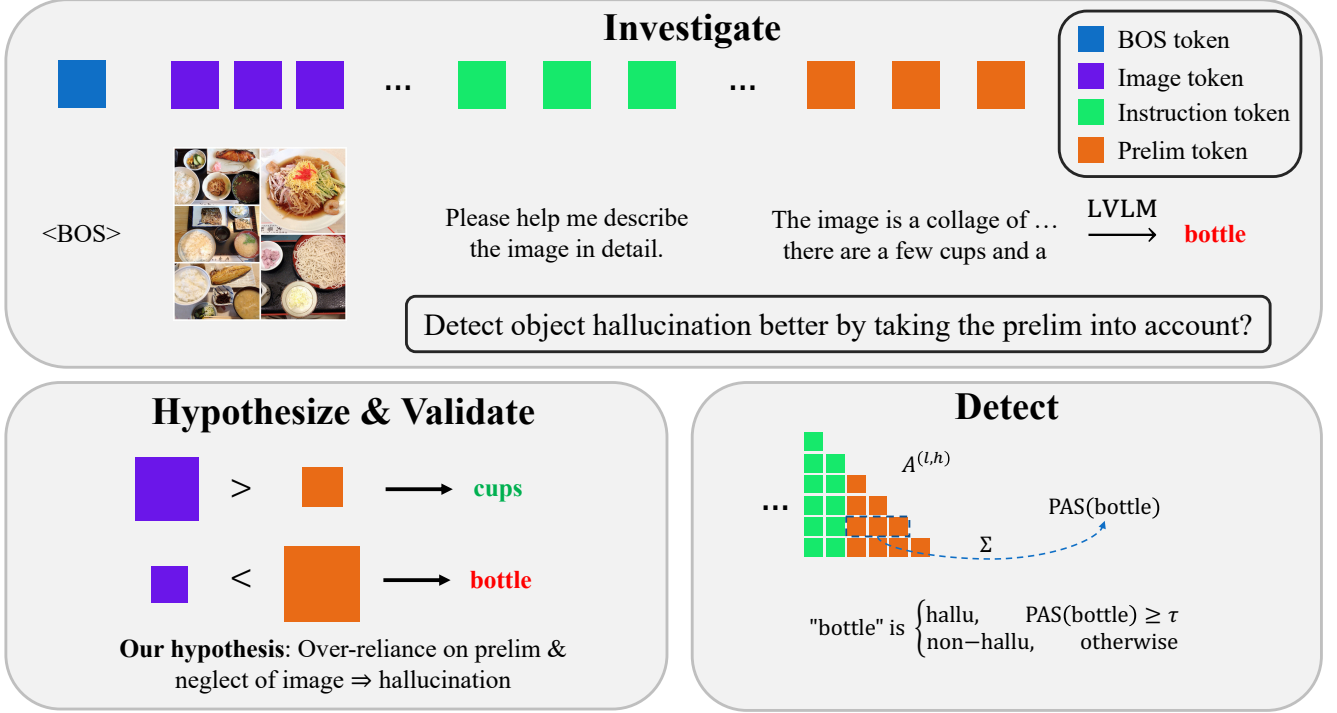
Figure 1. Illustration of our findings and proposed method. **Top (Investigate)**: We show that LVLM token generation depends on four token types and argue that prelim tokens are a vital, overlooked signal. **Bottom Left (Hypothesize & Validate)**: We visually illustrate our core hypothesis: hallucinations (e.g., "bottle") occur when the model over-relies on prelim tokens, while real objects (e.g., "cups") rely more on the image. **Bottom Right (Detect)**: Based on this insight, we propose PAS, which quantifies this over-reliance by summing the attention weights from prelim tokens to the object token.

which previous tokens significantly affect its generation. In this work, we leverage the attention weights to quantify how much prelim tokens affect generation of object tokens and link it to hallucination.

**Object hallucination detector.** When prompting the LVLM with an image $\mathbf{v}$, the output tokens $\mathbf{y}$ can contain references to objects, and we wish to determine if they are present or absence from the image $\mathbf{v}$ (corr. to real or hallucinatory). We define an *object hallucination scoring function* $D(y_k, \mathbf{y}, \mathbf{x})$ as a function that takes in an object mention $y_k \in \mathbf{y}$ and returns a real-valued score, where *higher* values indicate a *higher* chance of $y_k$ being a hallucination. Given this function, a detector can be obtained via thresholding with some threshold $\tau$, i.e., all $y_k$ with $D(y_k, \mathbf{y}, \mathbf{x}) \geq \tau$ is considered a hallucination.

## 4. Detecting Object Hallucination via Prelim

In this section, we first describe the LVLM token generating process. We then form our hypothesis on how prelim tokens are important in hallucination detection. Finally, we describe how to validate our hypothesis and discuss its implications, which leads to the methodology in Section 5.

### 4.1. Prelim matters for hallucination detection

As laid out in Section 3, the generation of a single token $y_k$ depends on image tokens $\mathbf{v}$, instruction tokens $\mathbf{t}$, and prelim tokens $\mathbf{y}_{<k}$. Crucially, the effect of each token type on the generation of $y_k$ depends on how much the model focus on them [12, 25], which can be quantified via the attention mechanism. Existing approaches [3, 17, 25] argue that in the normal operating mode, the model should mostly focus on the image for object generation, implying *low image attention → higher* chance of hallucination, and base on this to determine hallucinations. Importantly, they do not consider the prelim tokens for object hallucination detection. At first, it might be intuitive to disregard this source of information as the prelim contains much less information than the image, and is also less reliable because it is generated rather than given. However, in fact, the attention on the prelim is vital because if the model relies on unreliable information to generate tokens, then that behavior is undesired. Thus, we hypothesize that the inclusion of information from the prelim, which previous works missed, can help with hallucination detection. Our question is: "**Can we create a better object hallucination detector by taking the prelim into account?**"
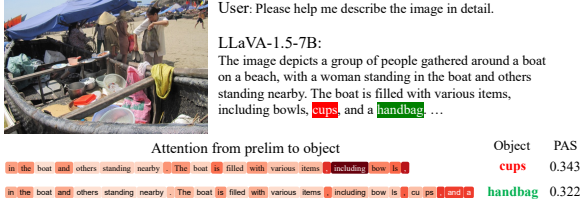
3

Figure 2. Visualization of Prelim Attention Score (PAS) for a single sample. We show the per-token attention for a suffix of the prelim for a hallucinated and a real object token. Darker red indicates higher attention to the corresponding object token, and a higher PAS score indicates a higher chance of hallucination.

In Fig. 3a, we visualize the sum of attention weights (formally defined in Eq. (10)) from prelim tokens to object tokens, categorized into real and hallucinatory tokens. This suggests that information from prelim tokens can be highly beneficial in determining hallucination. Specifically, the relation is *high prelim attention → higher* chance of hallucination, which is the inverse of that for image tokens. This raises the question: **when not attending to the image, why the model focuses on the prelim** instead of "certainly meaningless tokens" like the BOS [6]?"

We hypothesize that the model's focus on the prelim tokens reveals an *alternative operating mode*. In this mode, the model does not base its prediction on the image, possibly due to a lack of sufficient image-perceiving capabilities or because the image itself is noisy and is thus challenging to interpret. Regardless of the true reason, in this mode, the model predominantly depends on the prelim, but since this source of information is imperfect, it leads to more unreliable and hallucinatory predictions. If our hypothesis is true, then recognizing this mode could lead to a more effective hallucination detector, which motivates us to investigate it deeper. In the context of the object hallucination detection problem, we pose the following hypothesis (**H1**): "**If the generation of an object token depends heavily on the prelim tokens and not on the image tokens, then it is likely to be a hallucination.**"

To validate **H1**, we first define $Y_k$ as the random variable representing the token to be predicted at position $k$. **H1** says that there are two potential factors that determine $Y_k$: the prelim tokens $\mathbf{y}_{<k}$ and the image tokens $\mathbf{v}$. Moreover, it says that if $Y_k$ is not dependent on $\mathbf{v}$ then $y_k \sim Y_k$ is likely to be a hallucination. This notion can be precisely captured by the notion of mutual information (MI): if $I(\mathbf{v}; Y_k \mid \mathbf{y}_{<k}, \mathbf{t})$ is small then $y_k$ is likely to be a hallucination. Thus, we propose to use the MI as the detector:

$$D_{\mathrm{MI}}(y_k, \mathbf{y}, \mathbf{x}) = -I(\mathbf{v}; Y_k \mid \mathbf{y}_{<k}, \mathbf{t}), \qquad (3)$$

where

$$I(\mathbf{v}; Y_k \mid \mathbf{y}_{<k}, \mathbf{t}) = H(Y_k \mid \mathbf{y}_{<k}, \mathbf{t}) - H(Y_k \mid \mathbf{y}_{<k}, \mathbf{x}). \qquad (4)$$

The main challenge in computing Eq. (4) lies in obtaining the distribution of $Y_k \mid \mathbf{y}_{<k}, \mathbf{t}$, which is not directly available from the LVLM. Intuitively, this distribution can be thought of as the "prediction" given just the prelim, therefore, to compute $\Pr(Y_k = y \mid \mathbf{y}_{<k}, \mathbf{t})$, we propose to marginalize out the image $\mathbf{v}$ from Eq. (2):

$$\Pr(y \mid \mathbf{y}_{<k}, \mathbf{t}) = \mathbb{E}_{I \sim \mathcal{I}} \left[ \Pr \left( y \mid \mathbf{y}_{<k}, (I, \mathbf{t}) \right) \right]. \qquad (5)$$

We visualize the distribution of Eq. (3) for real and hallucinated tokens in Figure 3b and the detection AUROC in Table 1, which support our hypothesis **H1**. Further details on the implementation is described in Section 5.1.

### 4.2. Prelim Attention-based approach

A major issue with using the $D_{\mathrm{MI}}(y_k, \mathbf{y}, \mathbf{x})$ in Eq. (3) for detection is the estimation step (5). Averaging over some set of images $\mathcal{I}$ with size $|\mathcal{I}| = L$ necessitates $L$ extra forward passes of the model for each input $\mathbf{x}$. Thus, this approach demands $L + 1$ times more computation than the generation phase, which makes it highly impractical.

To overcome this, we propose an alternative detector based on the attention mechanism. From Fig. 3a, **the model assigns significantly higher attention weights to prelim tokens when generating hallucinatory object tokens, compared to real object tokens**. Let $s_{\mathrm{BOS}}, s_{\mathrm{img}}, s_{\mathrm{ins}}, s_{\mathrm{prel}}$ denote the sum of attention weights (formally defined in Eq. (10)) from respective token types to $y_k$. The softmax operation in self-attention layers requires that $s_{\mathrm{BOS}} + s_{\mathrm{img}} + s_{\mathrm{ins}} + s_{\mathrm{prel}} = 1$. This property, along with the fact that the attention mechanism is the only method of information exchange between tokens in a forward pass, suggests that attention scores can be used to quantify how much each token types affect generation of object tokens. Thus, in context of our hypothesis **H1**, we propose to directly use $s_{\mathrm{prel}}$ for hallucination detection:

$$D_{\mathrm{PAS}}(y_k, \mathbf{y}, \mathbf{x}) = s_{\mathrm{prel}}(y_k, \mathbf{y}, \mathbf{x}). \qquad (6)$$

In Section 5.2, we describe in details how to compute this score, and in Section 6.2, we show that this leads to an efficient detector.

## 5. Methodology

### 5.1. Mutual Information-based score

Since the estimation in Eq. (5) can introduce some errors, we also consider other candidates for computing $I(\mathbf{v}; Y_k \mid \mathbf{y}_{<k}, \mathbf{t})$, apart from the difference in entropy in Eq. (4). The detector in Eq. (3) can be generalized as follows:

| Variant | LLaVA-1.5-7B | | MiniGPT-4 | | Shikra | | Average |
|---------|--------|------------|--------|------------|--------|------------|---------|
| | MSCOCO | Pascal VOC | MSCOCO | Pascal VOC | MSCOCO | Pascal VOC | |
| Entropy diff (4) | 71.7 | 76.4 | 84.0 | 82.8 | 64.9 | 67.2 | 74.5 |
| KL div (8) | 73.8 | 81.1 | 86.5 | 86.2 | 73.1 | 80.4 | 80.2 |
| Logit diff (9) | 76.6 | 81.9 | **88.1** | **87.4** | 79.0 | 80.3 | 82.2 |
| PAS (6) | **84.2** | **85.1** | 85.6 | 85.4 | **84.5** | **85.3** | **85.0** |

Table 1. Comparing detection performance (AUROC, higher is better) of different choices of $\Delta$ in Eq. (7) against attention-based PAS across different models and datasets. All values are percentages, and best results are shown in **bold**.

$$D_\Delta(y_k, \mathbf{y}, \mathbf{x}) = \Delta\left(Y_k \mid \mathbf{y}_{<k}, \mathbf{t}; Y_k \mid \mathbf{y}_{<k}, \mathbf{x}\right) \quad (7)$$

in which $\Delta$ is a function that computes some "distance" between two discrete probability distributions. To realize Eq. (7), there are two points of focus: the estimation in Eq. (5), and the specification of $\Delta$.

**Estimation of** $\Pr(y \mid \mathbf{y}_{<k}, \mathbf{t})$**.** Since it is impossible to sample from the (unknown) distribution of all images, we estimate by averaging over a fixed set of images $\mathcal{I}$ instead. We lay out some strategies for a more stable estimation. We sample once per dataset, letting $\mathcal{I}$ be a subset of that dataset and use the same $\mathcal{I}$ for different prelims. Furthermore, since $y$ varies over the set of object tokens, ideally, we want $\mathcal{I}$ to contain all possible objects so that the model's true prediction tendency given a real object is observed for all object classes. This is achieved by randomly sampling an image that contains each class for all object classes. Thus, we have $|\mathcal{I}|$ = number of object classes in the dataset.

**Choice of** $\Delta$**.** First, we note that the entropy of a distribution is frequently difficult to compute. In our case, we can compute the entropy because the next token has only finite possibilities. Apart from the entropy difference as in (4), we consider two other candidates for quantifying the "distance" between the two discrete probability distributions. The first one is the KL Divergence, which is commonly used in the literature:

$$D_{\mathrm{div}}(y_k, \mathbf{y}, \mathbf{x}) = -D_{\mathrm{KL}}\left(Y_k \mid \mathbf{y}_{<k}, \mathbf{t} \| Y_k \mid \mathbf{y}_{<k}, \mathbf{x}\right). \quad (8)$$

Alternatively, we directly look at the difference in the predicted logit (pre-softmax) of the object token $y_k$. Intuitively, the logit can be interpreted as (the inverse of) the uncertainty for a prediction, which gives the following score:

$$D_{\mathrm{logit}}(y_k, \mathbf{y}, \mathbf{x}) = \left(\sum_{I \in \mathcal{I}} \frac{\Phi(y_k \mid \mathbf{y}_{<k}, (I, \mathbf{t}))}{|\mathcal{I}|}\right)$$
$$- \Phi(y_k \mid \mathbf{y}_{<k}, \mathbf{x}). \quad (9)$$

As shown in Table 1, the variants lead to functional detectors with varying degrees of effectiveness.

## 5.2. Prelim Attention Score

In this section, we describe in details the attention-based approach in Section 4.2. We start with some notations. Let $H$ be the number of attention heads in a layer of the LVLM. Let the attention weights at head $h$ of layer $l$ for the sequence $(\mathbf{x}, \mathbf{y})$ be $\mathbf{A}^{(l,h)} \in \mathbb{R}^{n \times n}$. Let $\mathbf{A}^{(l,h)}(k, j)$ denotes the attention from the $j$-th token in $(\mathbf{x}, \mathbf{y})$ to object token $y_k$. Then the prelim attention score (using layer $l$) is:

$$s_{\mathrm{prel}}(y_k, \mathbf{y}, \mathbf{x}) = \frac{1}{H} \sum_{h=1}^{H} \sum_{j=m+1}^{k-1} \mathbf{A}^{(l,h)}(k, j). \quad (10)$$

Different layers in the LVLMs might correspond to different stages of processing of the model [4, 14], which can lead to varying detection performance. Based on our ablation in Fig. 5 (discussed further in Sec. 6.3), we use the first layer (layer 0) as the default. Regarding multi-head attention, we average over all heads in a layer for simplicity, following [17]. While some work [30] suggests that selecting a certain subset of heads can be a better alternative than simply averaging over all heads, this adds additional complexity and requires validation so we leave it for future work. Furthermore, while this work focuses on the prelim, we also experiment with other scores for comparison. In particular, the scores $s_{\mathrm{BOS}}, s_{\mathrm{img}}, s_{\mathrm{ins}}$ can be defined similarly to Eq. (10), summing over the respective token types, and we experiment with them in Section 6.3.

As discussed earlier (Sec. 4.2), this attention-based approach is much more computationally efficient than the approach in Section 5.1 since it does not require extra passes and only relies on the attention weights, which are commonly computed during inference. Thus, in subsequent comparisons, unless otherwise stated, we use the prelim attention score of the first layer as our default method.

| Method | LLaVA-1.5-7B | | MiniGPT-4 | | Shikra | | Average |
|---|---|---|---|---|---|---|---|
| | MSCOCO | Pascal VOC | MSCOCO | Pascal VOC | MSCOCO | Pascal VOC | |
| NLL [37] | 56.5 | 64.0 | 62.1 | 73.0 | 54.3 | 63.1 | 62.2 |
| Entropy [26] | 71.7 | 64.3 | 69.8 | 62.9 | 71.4 | 64.4 | 67.4 |
| IC [16] | 75.1 | 64.6 | 76.4 | 67.7 | 76.0 | 71.3 | 71.9 |
| GLSim [27] | 64.1 | 69.4 | 63.6 | 62.0 | 67.8 | 66.6 | 65.6 |
| SVAR [17] | 81.5 | 82.9 | **88.0** | 84.5 | 71.9 | 72.9 | 80.3 |
| Ours | **84.2** | **85.1** | 85.6 | **85.4** | **84.5** | **85.3** | **85.0** |

Table 2. Object hallucination detection performance (AUROC, higher is better) on three models and across two datasets. For all models, greedy decoding with `max_new_tokens=512` is used. All values are percentages, and best results are shown in **bold**.

| Method | A | H | VRAM (GB) |
|---|---|---|---|
| Entropy | ✗ | ✗ | 16 |
| IC | ✗ | ✓ | 30 |
| GLSim | ✗ | ✓ | 19 |
| SVAR | ✓ | ✗ | 18 |
| Ours | ✓ | ✗ | 18 |

Table 3. Memory consumption of each detection method (**batch size 1**) when using LLaVA-1.5-7B in `float16` mode. Key: A (Attention) indicates the use of attention weights; H (Hidden states) indicates the use of hidden states. Note that the Entropy method uses neither Attention nor Hidden states, and thus represents the lowest memory consumption required.

| Method | Greedy | Beam search | Top-k | Nucleus |
|---|---|---|---|---|
| NLL | 56.5 | 55.7 | 58.3 | 59.9 |
| Entropy | 71.6 | 71.9 | 74.0 | 75.3 |
| IC | 75.1 | 73.0 | 76.7 | 77.9 |
| GLSim | 64.1 | 65.8 | 65.5 | 65.4 |
| SVAR | 81.5 | 81.3 | 79.9 | 80.4 |
| Ours | **84.2** | **84.0** | **83.5** | **84.0** |

Table 4. Effect of decoding methods on object hallucination detection performance (AUROC) for LLaVA-1.5-7B on MSCOCO. All values are percentages, and best results are shown in **bold**.

# 6. Experiments

## 6.1. Setup

**Models.** We conduct experiments on three popular LVLMs: LLaVA-1.5 [24], MiniGPT-4 [38], and Shikra [9]. All models have size 7B. For all models, greedy decoding with `max_new_tokens=512` is used. Additional results for bigger models and more implementation details are provided in the Appendix.

**Benchmarks.** We use the MSCOCO [23] dataset following previous works [17, 27], and additionally the Pascal VOC [11] dataset. For the MSCOCO dataset, following [27], we randomly sample 5,000 images from the `val2014` subset, and for the Pascal VOC dataset we use the full `val2012` subset, which contains 5,823 images. The MSCOCO and Pascal VOC dataset contain 80 and 20 object classes, respectively. These benchmarks contain, for each image, a list of objects known to be present in that image, and string matching against this list is used to discover object tokens, same as in CHAIR [28] hallucination evaluation. All models are prompted with "*Please help me describe the image in detail.*"

**Baselines.** We compare our detection method against other baselines, including *Logits*-based approaches: Negative Log-Likelihood (NLL) [37] and Entropy [26]; *Attention*-based approach: Summed Visual Attention Ratio (SVAR) [17]; and *Representation*-based approaches: Internal Confidence (IC) [16] and Global-Local similarity (GLSim) [27]. All baselines are evaluated on the same samples per dataset. For parameterized baselines, we use the parameters provided in the original papers for the models evaluated, if exists.

## 6.2. Results

In Table 1, we compare different realizations of Eq. (7) in terms of detection performance against the attention-based PAS. From this result, it can be seen that our information-theoretic formulation (Eq. (4)) leads to efficacious object hallucination detection, albeit with different performance depending on the choice of $\Delta$. We note with interest that our simple PAS method not only proves far more efficient (requiring no extra forward passes) but also achieves superior detection performance compared to all $D_\Delta$ variants. We hypothesize this is because PAS directly measures the model's internal information dependency for the current, specific input. In contrast, the $D_\Delta$ approach (Eq. (7)) is limited by the necessary but noisy estimation of the marginal probability distribution $\Pr(y_k \mid \mathbf{y}_{<k}, \mathbf{t})$ across a finite set of reference images $\mathcal{I}$. This suggests that direct attention flow is
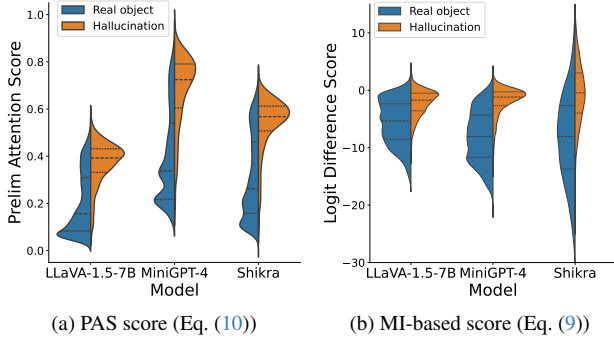
(a) PAS score (Eq. (10))  (b) MI-based score (Eq. (9))

Figure 3. Score distributions for real and hallucinated object tokens across different models on MSCOCO dataset. The dashed lines denote the quartiles for each distribution.
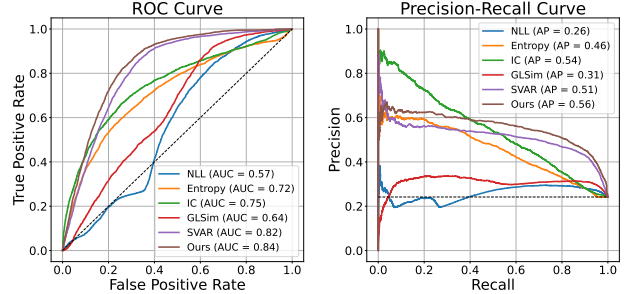


Figure 4. The ROC and PRC curves for object hallucination detection of our method and the baselines for LLaVA-1.5-7B on MSCOCO dataset. Dashed line indicates chance performance.

a cleaner and more robust signal of the model's "alternative operating mode" than distributional differences. Ultimately, the fact that both the probability-based $D_\Delta$ methods and the attention-based PAS yield effective detectors strongly supports our hypothesis **H1**: that prelim overdependence plays a significant role in object hallucination in LVLMs.

Figure 2 illustrates a qualitative example of our PAS method. In particular, we compare the attention from suffixes of the prelim to the two object tokens. The first object token receives more attention from the prelim, which indicates higher chance of hallucination. Notably, for the hallucinatory token, the attention seems to concentrate at a few tokens, while for the non-hallucinatory token it is more dispersed, which is potentially related to observations by [15].

Table 2 describes the main comparative study against other baselines on object hallucination detection. From the table, our PAS consistently outperforms the baselines across most models and datasets, highlighting the usefulness of prelim tokens attention for detection. Interestingly, the most competitive baseline is SVAR [17], which is also entirely attention-based, but they focus on image tokens instead. For a more detailed comparison, we plot the Receiver operating characteristic curve (ROC) and Precision-recall curve (PRC) of different methods for LLaVA-1.5-7B on MSCOCO in Figure 4. The curves suggest that our method remains effective irrespective of the choice of the threshold $\tau$. Additionally, Figure 3 depicts the distribution of our score for real and hallucinated tokens, showing that their distributions are effectively separated, as desired.

Finally, Table 3 indicates that our method incurs minimal GPU memory overhead, being on par with attention-based SVAR [17] and using 33% less *extra* memory per sample than hidden states-based GLSim [27], highlighting our method's efficiency.

### 6.3. Ablation studies

**Layer $l$ for computing attention score.** We perform an ablation test on the layer $l$ (Eq. (10)) chosen to compute PAS and visualize the results in Figure 5. Notably, the first layer tends to give the best results. Interestingly, this result aligns with existing research [4, 14] suggesting that earlier layers focus on gathering information and latter layers focus on processing that information. Based on this, we adopt the first layer (layer 0) to compute the attention scores for our method for the models evaluated in this paper.

**Attention from token types other than prelim.** Experimentally, we find that $s_{\text{BOS}}, s_{\text{ins}}, s_{\text{img}}$ have negative correlation with hallucination, while $s_{\text{prel}}$ has positive correlation with hallucination. Their detection performance is depicted in Table 5. We display both global (i.e., averaged over all layers) and layer 0 variants for a better overall picture.

A key question is whether PAS ($s_{\text{prel}}$) offers a genuine advantage over simply using image attention ($s_{\text{img}}$), which is the basis for methods like SVAR [17]. The softmax constraint in self-attention causes these scores to be correlated (as visualized in Figure 6). However, our results demonstrate they are not redundant signals. Table 5 clearly shows that $s_{\text{prel}}$ from layer 0 (84.8% avg) provides a more effective detection signal than $s_{\text{img}}$ from the same layer (82.1% avg). This suggests that while low image attention is a *symptom* of hallucination, high prelim attention is a more direct and potent *signal* of the model's shift to an unreliable "alternative operating mode".

Note that Table 5 also shows that some other attention scores (namely instruction tokens) exhibit significant detection capabilities, while not necessarily admitting a theoretically-inspired explanation similar to what we formalized for the prelim tokens. We theorize that this is because the softmax operation in the self-attention layers causes some degree of correlation between different attention scores. We visualize the pairwise correlation between the scores for two different models in Figure 6. Indeed,
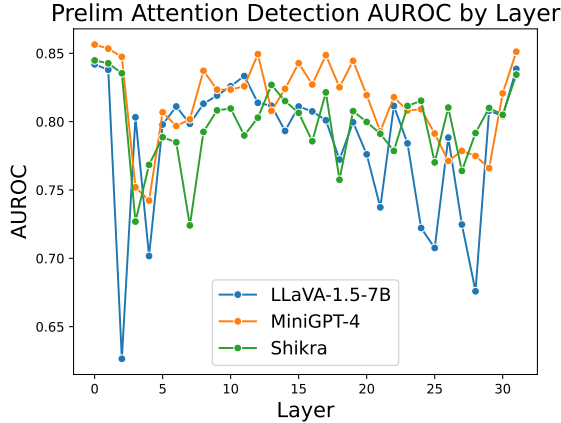
Figure 5. Comparison of different layers' performance in detecting object hallucination on the MSCOCO dataset.

for the models we considered, there is significant correlation between the prelim attention and instruction attention, which can explain why the instruction attention score performs well in Table 5.

**Effect of decoding strategies.** Decoding schemes can significantly affect hallucination rates of LVLMs [15]. Following existing works [17, 27], we use greedy decoding as the default evaluation setting due to its computational efficiency and simplicity. Here, we additionally investigate the effect of different decoding schemes on detection performance. We run the LVLMs with three additional distinct decoding strategies and perform the same quantitative experiments on object hallucination detection. We use $N_{\text{beams}} = 5$, $k = 10$, and $p = 0.9$ for beam search, top-k decoding, and nucleus decoding, respectively, and report the results in Table 4. Results suggest that our method consistently provides a significant gain in detection performance irrespective of the decoding scheme used, ensuring its practicality in real-world scenarios where decoding schemes can change depending on the task at hand [31].

## 7. Conclusion

We presented PAS, a training-free detector for object hallucinations in LVLMs that exploits a previously overlooked cue: *attention flow from preliminary, low-information output tokens to object tokens*. We support our method with an information-theoretic interpretation that links high prelim attention with an alternative unreliable operating mode of the LVLM. Experiments across models and datasets show that PAS achieves state-of-the-art object hallucination detection performance while being computationally efficient, requiring no extra passes. Our findings complement the literature on LVLM object hallucination, suggesting the exis-

| Method | LLaVA-1.5-7B | MiniGPT-4 | Shikra | Average |
|---|---|---|---|---|
| *Global attention scores* | | | | |
| Prelim | 83.7 | 86.2 | 84.4 | 84.8 |
| Instruction | 79.7 | 86.3 | 84.6 | 83.5 |
| Image | 75.0 | 88.7 | 72.4 | 78.7 |
| BOS | 74.0 | 74.9 | 55.5 | 68.1 |
| *Layer 0 attention scores* | | | | |
| Prelim | 84.2 | 85.6 | 84.5 | 84.8 |
| Instruction | 83.9 | 85.7 | 84.5 | 84.7 |
| Image | 84.1 | 85.3 | 77.0 | 82.1 |
| BOS | 83.1 | 85.4 | 82.1 | 83.5 |

Table 5. Hallucination detection results (AUROC) when using different types of attention scores on MSCOCO. *Global* score means averaging the attention weights over all heads in a layer and over all layers in a model.
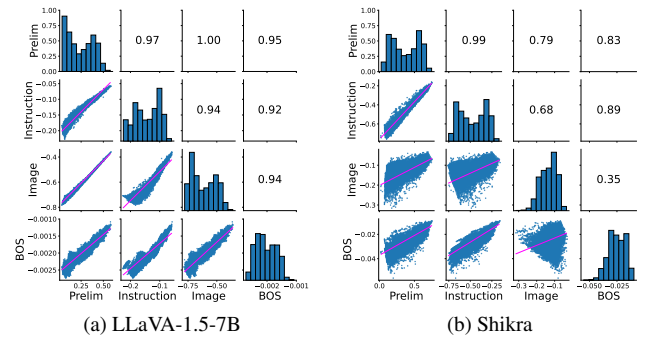


(a) LLaVA-1.5-7B (b) Shikra

Figure 6. Correlation between attention scores (layer 0) from different token types to object tokens for LLaVA-1.5-7B and Shikra on MSCOCO. Upper-right half shows Pearson correlation coefficients. All scores other than prelim attention score are negated so that high score indicates high likelihood of hallucination.

tence of an undesired operating mode of the model and its detection, thus paving the way for future development of more dependable and trustworthy LVLMs.

**Limitations and future work.** Our observations for the PAS score are based on the original self-attention proposed in [33] and is widely used by many LVLMs/LLMs. Besides full self-attention, other variants [2, 7] have been proposed to reduce the quadratic complexity of self-attention. For those variants, since the output token might no longer attend to all token types, our hypothesis/observations might not hold. Moreover, in this work, we focus on object existence hallucination; other types of hallucination such as relational/attribute hallucination are left for future work. Finally, our experiments focused on open-ended image description. Future work should investigate whether PAS remains an effective signal in other tasks, such as Visual Question Answering (VQA). It is plausible that in VQA, the varying questions introduce different prelim token dynamics, providing an interesting avenue for investigations.

## Acknowledgment

## References

[1] interpreting GPT: the logit lens — LessWrong. 2020. 2

[2] Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone, 2024. arXiv:2404.14219 [cs]. 8

[3] Wenbin An, Feng Tian, Sicong Leng, Jiahao Nie, Haonan Lin, QianYing Wang, Ping Chen, Xiaoqin Zhang, and Shijian Lu. Mitigating object hallucinations in large vision-language models with assembly of global and local attention. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 29915–29926, 2025. 2, 3

[4] Amit Ben Artzy and Roy Schwartz. Attend First, Consolidate Later: On the Importance of Attention in Different LLM

Layers. In *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 177–184, Miami, Florida, US, 2024. Association for Computational Linguistics. 5, 7

[5] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen Technical Report, 2023. arXiv:2309.16609 [cs]. 2

[6] Federico Barbero, Alvaro Arroyo, Xiangming Gu, Christos Perivolaropoulos, Petar Veličković, Razvan Pascanu, and Michael M. Bronstein. Why do LLMs attend to the first token? 2025. 4

[7] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The Long-Document Transformer, 2020. arXiv:2004.05150 [cs]. 8

[8] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. SpatialVLM: Endowing Vision-Language Models with Spatial Reasoning Capabilities. pages 14455–14465, 2024. 1

[9] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm's referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023. 6

[10] Xiang Chen, Chenxi Wang, Yida Xue, Ningyu Zhang, Xiaoyan Yang, Qiang Li, Yue Shen, Lei Liang, Jinjie Gu, and Huajun Chen. Unified Hallucination Detection for Multimodal Large Language Models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3235–3252, Bangkok, Thailand, 2024. Association for Computational Linguistics. 2

[11] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The Pascal Visual Object Classes (VOC) Challenge. *International Journal of Computer Vision*, 88(2):303–338, 2010. 6

[12] Xuan Gong, Tianshi Ming, Xinpeng Wang, and Zhihua Wei. DAMRO: Dive into the Attention Mechanism of LVLM to Reduce Object Hallucination. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7696–7712, Miami, Florida, USA, 2024. Association for Computational Linguistics. 3

[13] Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and Preventing Hallucinations in Large Vision Language Models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):18135–18143, 2024. 2

[14] Akshat Gupta, Jay Yeung, Gopala Anumanchipalli, and Anna Ivanova. How Do LLMs Use Their Depth?, 2025. arXiv:2510.18871. 5, 7

[15] Jieru Huang, Jinyue Yang, Guangzhi Ma, Danfeng Chen, Jason Kuen, Anima Anandkumar, and Gao Huang. Opera:

Over-trust penalty with rationale-guided decoding mitigates hallucination in vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13418–13427, 2024. 2, 7, 8

[16] Nicholas Jiang, Anish Kachinthaya, Suzanne Petryk, and Yossi Gandelsman. Interpreting and Editing Vision-Language Representations to Mitigate Hallucinations. 2024. 2, 6

[17] Zhangqi Jiang, Junkai Chen, Beier Zhu, Tingjin Luo, Yankun Shen, and Xu Yang. Devils in middle layers of large vision-language models: Interpreting, detecting and mitigating object hallucinations via attention lens. In *Proceedings of the computer vision and pattern recognition conference*, pages 25004–25014, 2025. 1, 2, 3, 5, 6, 7, 8

[18] Liqiang Jing, Ruosen Li, Yunmo Chen, and Xinya Du. Faith-Score: Fine-grained Evaluations of Hallucinations in Large Vision-Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 5042–5063, Miami, Florida, USA, 2024. Association for Computational Linguistics. 1, 2

[19] Sicong Leng, Jingjing Wang, Jinming Dong, Yang Gao, Weizhi Wang, Ee-Peng Lim, and Lidong Bing. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2

[20] Qing Li, Jiahui Geng, Chenyang Lyu, Derui Zhu, Maxim Panov, and Fakhri Karray. Reference-free Hallucination Detection for Large Vision-Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4542–4551, Miami, Florida, USA, 2024. Association for Computational Linguistics. 2

[21] Xin Li, Yunfei Wu, Xinghua Jiang, Zhihao Guo, Mingming Gong, Haoyu Cao, Yinsong Liu, Deqiang Jiang, and Xing Sun. Enhancing Visual Document Understanding with Contrastive Learning in Large Visual-Language Models. pages 15546–15555, 2024. 1

[22] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 292–305, Singapore, 2023. Association for Computational Linguistics. 1, 2

[23] Tsung-Yi Lin, Michael Maire, Serge J. 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*, pages 740–755. Springer, 2014. 6

[24] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual Instruction Tuning. In *Advances in Neural Information Processing Systems*, pages 34892–34916. Curran Associates, Inc., 2023. 6

[25] Shi Liu, Kecheng Zheng, and Wei Chen. Paying More Attention to Image: A Training-Free Method for Alleviating Hallucination in LVLMs. In *Computer Vision – ECCV 2024:*

*18th European Conference, Milan, Italy, September 29–October 4, 2024, Proceedings, Part LXXXIII*, pages 125–140, Berlin, Heidelberg, 2024. Springer-Verlag. 2, 3

[26] Andrey Malinin and Mark Gales. Uncertainty estimation in autoregressive structured prediction. In *International conference on learning representations*, 2021. 1, 2, 6

[27] Seongheon Park and Yixuan Li. GLSim: Detecting object hallucinations in LVLMs via global-local similarity. In *ICML 2025 workshop on reliable and responsible foundation models*, 2025. 1, 2, 6, 7, 8

[28] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4035–4045. Association for Computational Linguistics, 2018. 1, 2, 6

[29] Pranab Sahoo, Prabhash Meharia, Akash Ghosh, Sriparna Saha, Vinija Jain, and Aman Chadha. A Comprehensive Survey of Hallucination in Large Language, Image, Video and Audio Foundation Models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 11709–11724, Miami, Florida, USA, 2024. Association for Computational Linguistics. 1

[30] Sreetama Sarkar, Yue Che, Alex Gavin, Peter Anthony Beerel, and Souvik Kundu. Mitigating Hallucinations in Vision-Language Models through Image-Guided Head Suppression. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 12492–12511, Suzhou, China, 2025. Association for Computational Linguistics. 5

[31] Chufan Shi, Haoran Yang, Deng Cai, Zhisong Zhang, Yifan Wang, Yujiu Yang, and Wai Lam. A Thorough Examination of Decoding Methods in the Era of LLMs. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8601–8629, Miami, Florida, USA, 2024. Association for Computational Linguistics. 8

[32] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open Foundation and Fine-Tuned Chat Models, 2023. arXiv:2307.09288 [cs]. 2

[33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszko-reit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008, 2017. 8

[34] Wenyi Xiao, Ziwei Huang, Leilei Gan, Wanggui He, Haoyuan Li, Zhelun Yu, Fangxun Shu, Hao Jiang, and Linchao Zhu. Detecting and Mitigating Hallucination in Large Vision Language Models via Fine-Grained AI Feedback. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(24):25543–25551, 2025. 1, 2

[35] Dexuan Xu, Yanyuan Chen, Jieyi Wang, Yue Huang, Hanpin Wang, Zhi Jin, Hongxing Wang, Weihua Yue, Jing He, Hang Li, and Yu Huang. MLeVLM: Improve Multi-level Progressive Capabilities based on Multimodal Large Language Model for Medical Visual Question Answering. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 4977–4997, Bangkok, Thailand, 2024. Association for Computational Linguistics. 1

[36] Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. Woodpecker: hallucination correction for multimodal large language models. *Science China Information Sciences*, 67(12):220105, 2024. 1, 2

[37] Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and Huaxiu Yao. Analyzing and mitigating object hallucination in large vision-language models. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024. 1, 2, 6

[38] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. MiniGPT-4: Enhancing vision-language understanding with advanced large language models. In *The twelfth international conference on learning representations*, 2024. 6