# Mitigating Hallucinations in Vision-Language Models through Image-Guided Head Suppression

**Sreetama Sarkar**[1,*]    **Yue Che**[1,*]    **Alex Gavin**[3]    **Peter A. Beerel**[1]    **Souvik Kundu**[2]

[1]University of Southern California, Los Angeles, USA    [2]Intel Labs, USA
[3]Harvard-Westlake School, Los Angeles, USA

{sreetama,yueche,pabeerel}@usc.edu    souvikk.kundu@intel.com

## Abstract

Despite their remarkable progress in multimodal understanding tasks, large vision language models (LVLMs) often suffer from *"hallucination"*, generating texts misaligned with the visual context. Existing methods aimed at reducing hallucinations through inference time intervention incur a significant increase in latency. To mitigate this, we present **SPIN**, a task-agnostic attention-guided head suppression strategy that can be seamlessly integrated during inference ***without incurring any significant compute or latency overhead***. We investigate whether hallucination in LVLMs can be linked to specific model components. Our analysis suggests that hallucinations can be attributed to a dynamic subset of attention heads in each layer. Leveraging this insight, for each text query token, we selectively suppress attention heads that exhibit low attention to image tokens, keeping the top-$k$ attention heads intact. Extensive evaluations on visual question answering and image description tasks demonstrate the efficacy of SPIN in reducing hallucination scores up to $2.7\times$ while maintaining F1, and improving throughput by $1.8\times$ compared to existing alternatives. Code is available here.

## 1 Introduction

Large language models (LLMs) (Touvron et al., 2023) have revolutionized natural language understanding and generation, achieving state-of-the-art (SoTA) performance across numerous tasks. To extend these capabilities to vision-language tasks, large vision language models (LVLMs) (Liu et al., 2024b; Zhu et al., 2024) integrate an LLM backbone with vision encoders, mapping image inputs into the text embedding space. While this approach has enabled remarkable progress in vision-language understanding, LVLMs often generate output text misaligned with the visual context, a

---

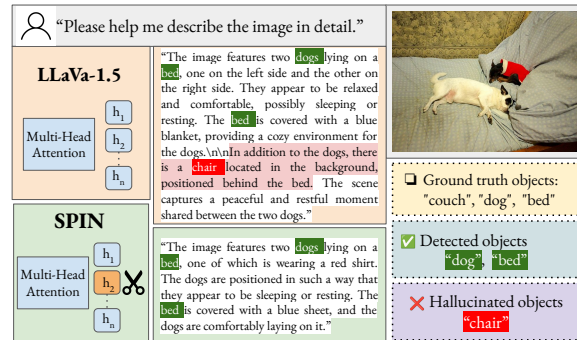*Equally contributing authors.



Figure 1: Caption generation using LLaVA-1.5 and SPIN. LLaVA-1.5's generated text description mentions a "chair" in the background, which is clearly a hallucinated object. SPIN mitigates hallucination while successfully identifying the objects present in the image.

phenomenon commonly referred to as "*hallucinations*", which undermines their reliability in critical domains such as healthcare, autonomous driving, and surveillance.

One key source of hallucination in LVLMs is textual bias, inherited from the pre-trained LLM backbone (Liu et al., 2024d; Chen et al., 2024a). To mitigate this, fine-tuning strategies based on Reinforcement Learning (RL) with human or AI-generated feedback have been proposed (Sun et al., 2024; Jing and Du, 2025). However, these methods are computationally intensive and often impractical for deployment in resource-constrained environments. Recent research has explored inference-only approaches, which modify the decoding pipeline (Huang et al., 2024) or apply contrastive decoding techniques (Leng et al., 2024; Liu et al., 2024d) to refine logit scores and reduce hallucinations. Although these methods are significantly more efficient compared to training-based approaches, they often struggle to sufficiently reduce hallucinations and suffer from significantly increased latency.

Besides, there is a lack of systematic understanding of how different model components contribute to hallucinations. In this paper, we investigate the role of attention heads in hallucinated text gen-

eration. We show that for each input token, a subset of attention heads in each layer disproportionately contributes to hallucinations. By identifying and suppressing these heads, we demonstrate that we can effectively reduce hallucinations in LVLMs while preserving SoTA model performance, as demonstrated in Figure 1.

**Our Contributions:** We perform a detailed analysis revealing that hallucination in LVLMs often stems from specific attention heads exhibiting insufficient attention to visual input. We characterize and quantify these "image-inattentive" heads across model layers. Based on this insight, we propose **SPIN**, **SuP**pressing image **IN**attentive heads, a novel, attention-guided head suppression strategy that can be seamlessly integrated during inference, irrespective of the decoding strategy or projection modules. SPIN offers a highly efficient solution since it achieves substantial hallucination reduction without requiring model retraining, and crucially, introduces ***no additional computational overhead or latency during inference***. We evaluate our approach on visual question answering (VQA) and image caption generation tasks. SPIN reduces hallucination scores up to $2.7\times$ over existing methods, while improving throughput by up to $1.8\times$.

## 2 Background and Related Work

### 2.1 LVLM Preliminaries

LVLMs typically comprise four key components: a text tokenizer, an image encoder, a projector, and a language decoder. The text tokenizer processes the language input by segmenting it into discrete tokens and converting them into text embeddings, henceforth referred to as *text tokens*. Similarly, the image encoder partitions the input image into patches and transforms them into corresponding visual embeddings. These visual embeddings are then mapped into the text embedding space through the projector, producing *vision tokens*. Finally, the vision and text tokens are concatenated and fed into the language decoder.

The language decoder consists of a series of transformer encoder blocks, each consisting of a multi-head attention (MHA) layer followed by a feed-forward network (FFN). The MHA takes in $N$ input tokens, consisting of both text and vision tokens, each of embedding dimension $d$, $X \in \mathbb{R}^{N \times d}$ and maps them into Query, Key and Value matrices $(Q, K, V) \in \mathbb{R}^{N \times d}$. MHA then computes low-dimensional projections of $(Q, K, V)$ given

by $(Q^i, K^i, V^i) \in \mathbb{R}^{N \times d_k}$ for each head $i$ where $d_k = d/H$, $H$ denoting the number of attention heads. The scaled dot product attention for each head is then computed using

$$h_i = \text{Softmax}\left(\frac{Q^i K^{iT}}{\sqrt{d_k}}\right) V^i \qquad (1)$$

The results for each head are concatenated and a final projection matrix $W_o \in \mathbb{R}^{d \times d}$ is applied to obtain the MHA output.

$$\text{MHA}_{Q,K,V} = \left(\bigoplus_{i=1}^{H} h_i\right) W_o \qquad (2)$$

The output text sequence is generated based on a *decoding strategy*, which determines how the next token is selected from the probability distribution obtained from the language decoder. Some of the common approaches include *greedy search*, which deterministically selects the token with the highest probability at each step; *beam-search* (Freitag and Al-Onaizan, 2017), which explores multiple candidate token sequences ("beams") simultaneously to find a higher-scoring overall sequence; and *nucleus sampling* (Holtzman et al., 2020), which samples the next token randomly from the smallest set of tokens whose cumulative probability exceeds a given threshold. Additionally, a repetition penalty is often applied during decoding to discourage the model from generating repetitive text by reducing the likelihood of tokens that have recently appeared in the sequence.

### 2.2 LVLM Hallucination and Mitigation

Hallucinations can manifest as factual errors to a user's query or inaccurate image description resulting in non-existent objects, incorrect object attributes or relationships (Liu et al., 2024a). Several causes of hallucinations have been identified including biased training data (Liu et al., 2023), the inability of vision encoders to accurately ground images (Jain et al., 2024), misalignment among different modalities (Liu et al., 2024b), and insufficient context attention in LLM decoders (Huang et al., 2024; Liu et al., 2024d). Existing hallucination mitigation approaches can be broadly classified into *training-based* and *training-free* methods.

**Training-based Mitigation:** (Sun et al., 2024) adapt the Reinforcement Learning from Human Feedback (RLHF), originally developed for text-only models, to the vision-language setting, training VLMs to maximize simulated human rewards.

They propose Factually-Augmented RLHF, where the reward model is given access to additional ground truth information, such as image captions, to improve its assessment of factual correctness. In contrast, FGAIF (Jing and Du, 2025) replaces human supervision with fine-grained AI-generated feedback. It segments model responses into sub-sentences and uses AI models to detect hallucinations relating to object existence, attributes, and relationships. LACING (Zhao et al., 2024) addresses language bias in LVLMs by introducing a multimodal dual-attention mechanism and soft-image guidance. It constructs separate attention streams for visual and textual inputs to improve grounding and alignment. Despite their effectiveness, training-based methods are computationally expensive, requiring large-scale resources (e.g., 8× A100 GPUs with 40GB memory each), making them impractical in resource-constrained settings.

**Training-free Mitigation:** Recently, training-free approaches have been proposed to mitigate the issue of visual context neglect in LVLMs. OPERA(Huang et al., 2024) observes that hallucination arises from generating new tokens based on limited summary tokens in which attention aggregation occurs, leading the model to ignore the image context. To address this, OPERA introduces a beam-search variant with a weighted scoring mechanism that downranks candidate sequences exhibiting over-trust patterns. However, its applicability is restricted to beam-search, which is significantly slower than greedy decoding. Other techniques leverage contrastive decoding, which requires multiple forward passes, increasing latency. VCD (Leng et al., 2024) builds on the observation that increased visual uncertainty drives models to rely on language priors, thereby amplifying hallucinations. It contrasts the output distributions conditioned on original and distorted visual inputs to identify and suppress hallucinated content. PAI (Liu et al., 2024d) mitigates the reduced attention to image tokens through a two-step approach: (1) amplifying attention to image tokens and (2) refining logit scores by subtracting logits computed without image prior to eliminate the text bias. DAMRO (Gong et al., 2024) attributes high-attention outlier tokens scattered in the background of the image as the cause for hallucinations and proposes to mitigate the influence of outlier tokens using contrastive decoding. HALC (Chen et al., 2024b) is based on contrasting distributions with
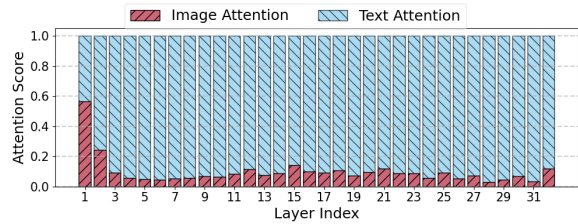


Figure 2: Average attention allocated by the current token to the preceding vision and text tokens in CHAIR. Image tokens receive <10% of total attention from layer 3, while constituting ∼76-92% of the input.

different visual contexts and using visual matching scores for candidate selection in beam-search. While all the above approaches introduce latency overhead, HALC is reported to cost around 2.4× of the normal greedy decoding time. Recently, (Liu et al., 2025) introduced a test-time latent space steering approach through visual and textual intervention (VTI) using pre-computed steering vectors, whereas ICT (Chen et al., 2025) proposes image and object-level intervention to apply targeted activation shifts to selected attention heads identified through binary classifiers, requiring additional parameters. In contrast, we propose a training-free mitigation strategy that introduces no additional overhead by identifying the role of attention heads in hallucination.

## 2.3 Role of Attention Heads

Earlier works (Voita et al., 2019) analyzed the role of individual attention heads in neural machine translation, showing that most heads can be pruned without significantly affecting model performance. More recently, routing to expert attention heads (Jin et al., 2024; Zheng et al., 2025) was introduced for LLMs and vision transformers, allowing each input token to dynamically select appropriate heads, thereby improving performance. (Zhou et al., 2024) further highlighted the impact of specific attention heads on LLM safety, demonstrating that ablating a small subset of heads increases the attack success rate. Recently, (Kundu et al., 2025) explored the impact of LLM weight (you, 2024) and KV quantization (Kang et al., 2024) on accuracy as well as hallucination performance of LLMs. However, the role of attention heads in the context of hallucination remains unexplored.

## 3 Methodology

**Motivational Study:** LVLMs struggle with inefficient image attention (Chen et al., 2024a; Liu et al., 2024d). In Figure 2, we illustrate how the current

token distributes attention across preceding vision and text tokens (constituting system prompt, instruction prompt, and generated output) averaged over all output tokens in the image caption. Our analysis reveals that in deeper layers, generated text tokens allocate $<10\%$ of their attention to vision tokens, despite vision tokens comprising $\sim$76-92% of the input. This imbalance causes the model to ignore the context provided by the visual input (taken as the "fact"), potentially leading to hallucinations. To mitigate this, for each query text token, we identify a fraction of heads that allocates the least cumulated attention to vision tokens. We then present a training-free strategy to dynamically suppress attention heads, reducing the imbalance in attention and thereby enhancing model performance.

**SPIN Multi-Head Attention:** We propose to mitigate problematic attention heads across layers using a dynamic mask $m_i$ for each attention head $i$. $m_i$ is obtained based on the attention of the current text query token $q^i \in \mathbb{R}^{1 \times d_k}$ to key vision tokens, having length $N_v$ given by:

$$\mathbf{A}_v = \mathbf{A}_{tot}[I_{\text{start}} : I_{\text{end}}], \quad \mathbf{A}_{tot} = q^i K^{iT} \quad (3)$$

Here, $\mathbf{A}_{tot} \in \mathbb{R}^{1 \times N}$ denotes the attention of $q^i$ to $N$ input tokens, $K^i \in \mathbb{R}^{N \times d_k}$ represents the key matrix for $N$ tokens, and $\mathbf{A}_v \in \mathbb{R}^{1 \times N_v}$ denotes the attention score for $q^i$ to the vision tokens only. $I_{\text{start}}$ and $I_{\text{end}}$ denotes the start and end indices of vision tokens. The mask $m_i$ is defined as:

$$m_i = \begin{cases} 1 & \text{if } i \in \text{top-}k(\sum_{j=1}^{N_v} A_v[j]) \\ \alpha & \text{otherwise} \end{cases}$$

Thus, $m_i$ is set to 1, identifying that we should keep the head intact if the $i^{th}$ attention head belongs to the top $k$ highest values across $n$ heads. Otherwise, $m_i$ is set to $\alpha$, which is the suppression factor used to reduce the influence of the $i^{th}$ head. If $\alpha = 0$, the head is completely suppressed or effectively pruned. We denote the ratio of suppressed heads as $r$ where $r = (1 - k/H)$. The choice of layers for suppressed heads, suppressed head ratio $r$ and suppression factor $\alpha$ are the key hyperparameters in our approach. The final multi-head attention in SPIN is computed as:

$$\text{MHA}_{Q,K,V,m} = \left( \bigoplus_{i=1}^{H} (m_i \cdot h_i) \right) W_o. \quad (4)$$

**Head and Suppression Factor Selection:** To identify the dynamic subset of problematic attention

heads for a specific model and task, we adopt an efficient systematic three-stage approach. *First,* we vary $r$ to find the value that achieves the best reduction in hallucinations without a significant drop ($\sim$3%) in F1. Here, we prune attention heads across all layers equally setting $\alpha = 0$. Adjusting these parameters is explored next. *Second,* our analysis (Section 4.8) reveals that attention heads in earlier layers tend to contribute more to hallucinations than those in later layers. For specific models and tasks, we find the optimal number of early layers in which we should prune attention heads to minimize hallucination scores. *Finally,* we explore increasing $\alpha$ for the selected heads to mitigate the drop in F1 observed in step 1. In particular, a higher $\alpha$ results in improving F1 score at the cost of increasing hallucinations. We choose an $\alpha$ that provides an optimal trade-off between hallucinations and F1.

# 4 Experimental Results

## 4.1 Experimental Setup

We evaluate our approach on LLaVA-1.5 (7B, 13B) (Liu et al., 2024b), LLaVA-Next (Liu et al., 2024c), MiniGPT-4 (Zhu et al., 2024), Shikra (Chen et al., 2023), Qwen-VL (Bai et al., 2023), and Qwen2.5-VL (Bai et al., 2025) models using POPE (Li et al., 2023), CHAIR (Rohrbach et al., 2018), MMHal-Bench (Sun et al., 2024), MME (Fu et al., 2023), MMMU (Yue et al., 2024), and GPT-4o assisted evaluation. While POPE and CHAIR evaluate object hallucinations, with POPE using Yes-No questions and CHAIR assessing image descriptions, MMHal-Bench takes into account logical considerations like object count, attributes, and relationships. We report results using different decoding strategies including greedy, nucleus sampling, and beam-search (with beam width of 5), and compare our approach with existing methods like OPERA, VCD, PAI, and DAMRO for each decoding strategy. Notably, OPERA is only compatible with beam-search, while VCD is specifically designed for nucleus sampling. We report DAMRO results only for LLaVA models, as its applicability to more complex projection modules like Q-Former remains unclear, a limitation also acknowledged in the original paper. The hyperparameters used for the existing methods for each model are taken from their respective papers (further details are provided in the Appendix A.1). Additionally, to evaluate the efficiency of our method, we report throughput mea-
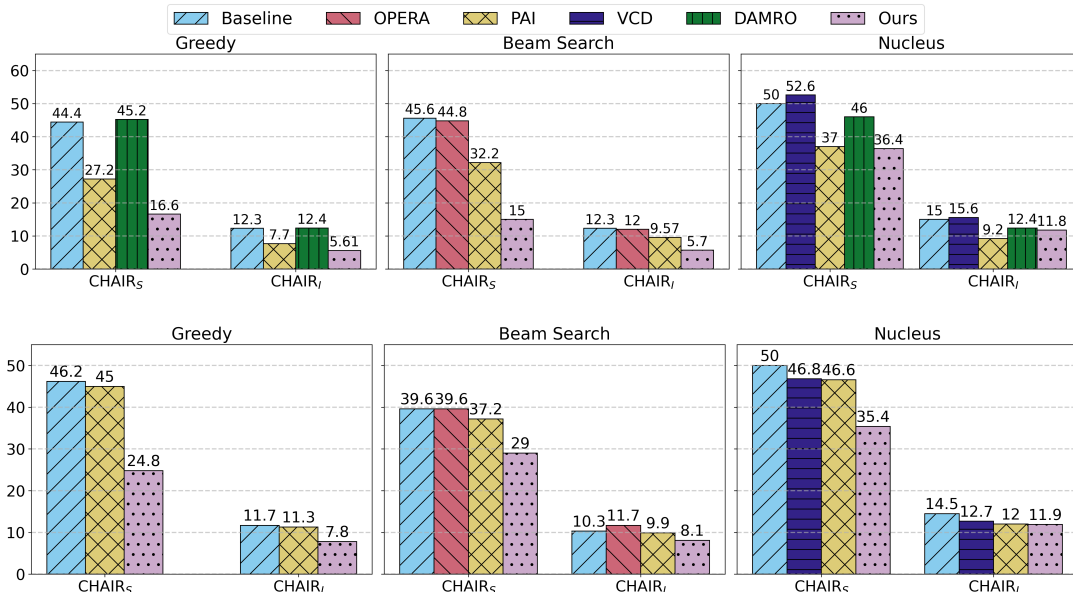
Figure 3: CHAIR scores for SPIN compared with existing approaches for greedy, beam-search, and nucleus sampling based decoding on LLaVA-1.5 (7B) (top) and Qwen-VL (bottom).

| Model | Method | Layers | $r$ | $\alpha$ | $C_S$ | $C_I$ | F1 |
|-------|--------|--------|-----|----------|-------|-------|-----|
| LLaVA-1.5 (7B) | Baseline | - | - | - | 44.4 | 12.3 | **77.8** |
| | PAI | - | - | - | 27.2 | 7.7 | 76.8 |
| | DAMRO | - | - | - | 45.2 | 12.4 | **77.8** |
| | VTI | - | - | - | 35.8 | 11.1 | 76.8 |
| | SPIN | 1∼32 | 0.05 | 0.08 | 26.4 | 7.6 | 77.6 |
| | SPIN | 1∼32 | 0.05 | 0.01 | **16.6** | **5.6** | 74.6 |
| LLaVA-1.5 (13B) | Baseline | - | - | - | 41.4 | 10.9 | 78.9 |
| | PAI | - | - | - | 37.4 | 9.2 | 79.2 |
| | DAMRO | - | - | - | 41.2 | 11.0 | 78.7 |
| | SPIN | 1∼16 | 0.10 | 0.0 | 30.6 | 8.3 | **79.6** |
| | SPIN | 1∼20 | 0.10 | 0.0 | **29.2** | **7.9** | 79.1 |
| MiniGPT-4 | Baseline | - | - | - | 31.4 | 11.1 | **70.6** |
| | PAI | - | - | - | 19.8 | 8.4 | 69.7 |
| | SPIN | 1∼16 | 0.18 | 0.0 | 21.0 | **6.2** | 68.8 |
| | SPIN | 1∼16 | 0.18 | 0.05 | **17.6** | 8.4 | 68.4 |
| Qwen-VL | Baseline | - | - | - | 46.2 | 11.7 | 76.5 |
| | PAI | - | - | - | 45.0 | 11.3 | 76.6 |
| | SPIN | 1∼20 | 0.30 | 0.08 | 29.6 | 8.5 | **78.0** |
| | SPIN | 1∼20 | 0.30 | 0.001 | **24.8** | **7.8** | 76.7 |
| Shikra | Baseline | - | - | - | 55.2 | 14.0 | 75.39 |
| | PAI | - | - | - | 34.0 | 8.5 | **76.4** |
| | SPIN | 1∼32 | 0.40 | 0.0 | 30.6 | 8.0 | 75.41 |
| | SPIN | 1∼32 | 0.45 | 0.0 | **24.4** | **7.1** | 74.3 |

Table 1: CHAIR evaluation with Greedy decoding.

sured in tokens per second. All experiments are run in Pytorch using NVIDIA RTX A6000 GPUs.

## 4.2 CHAIR Evaluation

**Method:** The Caption Hallucination Assessment with Image Relevance (CHAIR) (Rohrbach et al., 2018) is a widely used metric for evaluating object hallucinations in image captioning tasks. It features two variants: the per-instance metric $C_i$, which indicates the fraction of object instances that are hallucinated, and the per-sentence metric $C_s$,

which indicates the fraction of sentences containing a hallucinated object.

CHAIR compiles a set of ground-truth objects for each image. Any object included in the image caption that is not present in the ground-truth object set is classified as a "hallucinated object". We evaluate our approach on randomly sampled 500 images from the COCO 2014 (Lin et al., 2014) validation set with the prompt *"Please help me describe the image in detail."*, following the same setup as used in Liu et al. (2024d); Huang et al. (2024). We report F1 scores alongside CHAIR scores to ensure that reducing hallucinations does not come at the expense of missing correct objects.

**Evaluation on greedy decoding:** We present CHAIR results for Baseline, PAI, DAMRO, VTI, and SPIN using Greedy decoding in Table 1. A lower CHAIR score indicates fewer hallucinated objects and a higher F1 score indicates higher caption accuracy. SPIN reduces $C_i$ and $C_s$ by 2.2× and 2.7× respectively over baseline LLaVA 7B with ∼3% degradation in F1 when 5% heads are suppressed uniformly across all layers. The drop in F1 can be mitigated using a higher scaling factor while still maintaining lower CHAIR scores compared to existing methods. For LLaVA 13B, Qwen-VL and Shikra, SPIN outperforms baseline even in terms of F1 scores while reducing $C_i$ and $C_s$ by 1.38× and 1.42× for LLaVA 13B, 1.5× and 1.9× for Qwen-VL, and 1.97× and 2.3× for Shikra. For MiniGPT-4, SPIN reduces $C_i$ and $C_s$ by 1.8× and

| Model | Mode | Layers | Supp. Heads | Scale Factor | Random | | Popular | | Adversarial | | Overall | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 |
| LLaVA-1.5 (7B) | Baseline | - | - | - | 86.77 | 85.21 | 85.73 | 84.36 | 84.57 | 83.27 | 85.69 | 84.28 |
| | PAI | - | - | - | 88.53 | 87.54 | 87.40 | 86.53 | 85.43 | 84.59 | 87.12 | 86.22 |
| | DAMRO | - | - | - | 86.73 | 85.18 | 85.77 | 84.39 | 84.53 | 83.22 | 85.68 | 84.26 |
| | SPIN | 1∼32 | 0.20 | 0.1 | **89.47** | **88.66** | **88.53** | **87.84** | **86.63** | **85.89** | **88.21** | **87.46** |
| LLaVA-1.5 (13B) | Baseline | - | - | - | 86.77 | 85.81 | 86.27 | 85.56 | 82.33 | 82.07 | 85.12 | 84.48 |
| | PAI | - | - | - | 87.67 | 87.03 | 86.97 | 86.60 | 83.03 | 83.14 | 85.89 | 85.59 |
| | SPIN | 1∼20 | 0.35 | 0.0 | **91.17** | **90.65** | **88.83** | **88.46** | **85.53** | **85.49** | **88.51** | **88.20** |
| MiniGPT-4 | Baseline | - | - | - | 83.35 | 81.28 | 75.23 | 73.79 | 75.66 | 73.52 | 78.08 | 76.20 |
| | PAI | - | - | - | 85.25 | 83.62 | **75.97** | 74.93 | **77.91** | 76.10 | **79.71** | 78.22 |
| | SPIN | 1∼24 | 0.05 | 0.0 | **85.64** | **84.56** | 75.36 | **75.34** | 77.50 | **76.86** | 79.50 | **78.92** |
| Shikra | Baseline | - | - | - | 80.77 | 80.73 | 78.06 | 78.71 | 75.99 | 77.22 | 78.27 | 78.89 |
| | PAI | - | - | - | 80.73 | 80.30 | 77.90 | 78.52 | 75.90 | 77.08 | 78.18 | 78.63 |
| | SPIN | 1∼16 | 0.20 | 0.001 | **82.04** | **81.56** | **79.15** | **79.38** | **76.78** | **77.56** | **79.32** | **79.50** |

Table 2: Multi-turn POPE evaluation across random, popular, and adversarial splits using Greedy Decoding.

1.5× with 1.8% degradation in F1. We also outperform PAI, DAMRO, and VTI for all models using greedy decoding. These significant improvements in the CHAIR scores quantify the benefits of the head pruning. We provide additional CHAIR evaluation results for LLaVA-Next and Qwen2.5-VL in the Appendix Table 12.

In general, we observe that the problematic heads lie in the first 16 to 20 layers, whereas for some models like LLaVA 7B and Shikra, they are uniformly distributed across all layers. The fraction of suppressed heads ($r$) varies significantly across models, ranging from 5% in LLaVA 7B to 45% in Shikra. Shikra's higher head suppression ratio likely indicates its attention heads perform more distributed or redundant functions, as opposed to LLaVA, where heads are more specialized. A detailed ablation study analyzing the impact of layer position, the fraction of suppressed heads, and the choice of scaling factor is provided in Section 4.8.

**Evaluation on other decoding strategies:** In Figure 3, we compare SPIN to existing methods for different decoding modes for LLaVA 7B and Qwen-VL. We compare with OPERA for beam-search decoding, VCD for nucleus sampling, DAMRO for greedy and nucleus sampling, and PAI for all three modes. Our results demonstrate that we achieve up to 3× lower CHAIR scores compared to existing methods for greedy and beam-search. For nucleus sampling, SPIN exhibits a slight degradation in performance likely due to the variability in token selection, as further elaborated in Section 4.3.

**Mitigating Repetition from Head Pruning:** In our experiments, we observed that SPIN occasionally leads to repetitive outputs for MiniGPT-4. While SPIN effectively reduces hallucinations and

| Method | Rep. Penalty | Scale Factor | CHAIR$_S$ (↓) | CHAIR$_I$ (↓) | F1 (↑) | Caption Length |
|---|---|---|---|---|---|---|
| Baseline | 1.0 | - | 31.4 | 11.1 | 70.6 | 84.1 |
| SPIN (r=0.15, Layers 1∼24) | 1.0 | 0.0 | 19.6 | 6.8 | 68.2 | 321.3 |
| | 1.05 | 0.0 | 19.2 | 7.7 | 68.9 | 100.1 |
| | 1.08 | 0.0 | 20.2 | 7.8 | 68.1 | 57.4 |
| | **1.1** | **0.0** | **19.2** | **7.8** | **67.5** | **42.9** |
| SPIN (r=0.18, Layers 1∼16) | 1.0 | 0.0 | 21.0 | 6.2 | 68.8 | 271.6 |
| | 1.1 | 0.0 | 17.4 | 8.3 | 67.7 | 43.0 |
| | **1.1** | **0.05** | **17.6** | **8.4** | **68.4** | **44.7** |
| | 1.1 | 0.08 | 21.6 | 9.2 | 69.4 | 46.2 |

Table 3: Analysis using repetition penalty on MiniGPT-4

maintains F1 scores close to the baseline, it increases the length of generated sequences by up to 3.8×. To mitigate this, we leverage the repetition penalty during decoding, as shown in Table 3. While a high value can negatively impact F1, a repetition penalty of 1.1 alleviates repetitions without compromising overall output quality, as evidenced by caption length and F1 performance.

### 4.3 POPE Evaluation

**Method:** Polling-based Object Probing Evaluation (POPE) (Li et al., 2023) is used for assessing hallucinations in VQA tasks by querying models with "Is there a <object> in the image?". The objects are drawn from three splits: *random* (any dataset object), *popular* (most frequent objects), and *adversarial* (closely related but misleading objects). We evaluate on 500 randomly selected COCO 2014 (Lin et al., 2014) validation images, with six questions per image from each split. We present results for multi-turn POPE evaluation (Liu et al., 2024d), where earlier responses are appended to the input context, increasing the context length and amplifying the chances of image neglect.
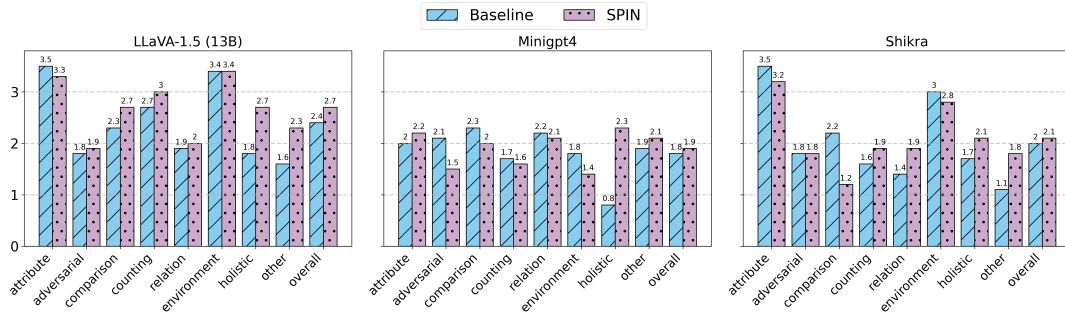
Figure 4: MMHal-Bench evaluation on LLaVA-1.5 (13B), MiniGPT-4 and Shikra.

| Decoding | Mode | Accuracy | F1 |
|---|---|---|---|
| Beam Search | Baseline | 83.58 | 81.14 |
| | OPERA | 83.67 | 81.29 |
| | PAI | 84.52 | 82.80 |
| | SPIN | **86.15** | **84.93** |
| Nucleus | Baseline | 82.97 | 81.70 |
| | VCD | 84.42 | 83.45 |
| | PAI | 84.96 | 83.99 |
| | DAMRO | **85.69** | **84.28** |
| | SPIN | 82.49 | 81.83 |

Table 4: Multi-turn POPE evaluation using Beam Search and Nucleus Sampling decoding for LLaVA-1.5 (7B).

**Evaluation on greedy decoding:** Table 2 reports POPE results for LLaVA (7B and 13B), MiniGPT-4 and Shikra across all three splits using greedy decoding. SPIN consistently improves accuracy and F1 across all three splits for LLaVA and Shikra models, achieving gains of up to 3.4% in accuracy and 3.7% in F1 over the baseline. DAMRO performs similarly or worse than baseline for greedy decoding. SPIN consistently outperforms PAI in both accuracy and F1 for LLaVA and Shikra, achieving up to a 2.6% gain in both metrics. For MiniGPT-4, SPIN attains the highest F1 score, while PAI achieves slightly better accuracy on the popular and adversarial splits. We provide POPE evaluation results for LLaVA-Next and Qwen2.5-VL in the Appendix Table 13.

**Evaluation on other decoding strategies:** In Table 4, we report overall accuracy and F1 scores for multi-turn POPE evaluation using beam-search and nucleus sampling decoding. We observe that SPIN yields the best results for beam-search, whereas DAMRO performs the best for nucleus sampling. We observe that SPIN performs more effectively with greedy decoding and beam-search, but is less effective under nucleus sampling. This may be attributed to the inherent randomness in sampling-based methods, which introduce greater variability in token selection. As a result, even after suppressing hallucination-prone heads, the stochastic nature

of nucleus sampling can still lead to hallucinated outputs, limiting the impact of structural interventions like head suppression.

### 4.4 MMHal-Bench Evaluation

**Method:** To extend our analysis to more complex and logically challenging benchmarks, we evaluate using MMHal-Bench (Sun et al., 2024). This benchmark comprises 96 carefully curated image-question pairs based on images from OpenImages (Krasin et al., 2016), each accompanied by ground-truth answers. The pairs span 12 common object meta-categories derived from COCO, and the questions test nuanced reasoning across eight categories: object attributes, adversarial objects, comparisons, counting, spatial relations, environment, holistic descriptions, and others—which includes cases where models fail to recognize text, misinterpret icons, or incorrectly reason about the observed visual content. We generate responses using both the baseline model and SPIN, and then employ GPT-4 to score each response based on its agreement with the ground-truth answers.

**Evaluation Results:** The final evaluation report including per-category scores for each model, with the overall performance measured as the average across all categories, is presented in Figure 4. We observe that SPIN improves overall scores over the baseline for all three models: LLaVA-1.5, MiniGPT-4 and Shikra, with the maximum improvements observed in *holistic description* and *other* categories. Note, we use the same values for Layers, $r$, and $\alpha$ as reported for CHAIR evaluation.

### 4.5 GPT-4o Assisted Evaluation

We use GPT-4o to compare the responses of baseline and SPIN models using Greedy decoding on a scale of 1 to 10 using two metrics: (1) *Accuracy (A):* measuring if the description is accurate with respect to the image content, and (2) *Detailedness*

| Model | Method | Perception | | | | | | | | | | | Cognition | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Exis-tence | Count | Posi-tion | Color | Posters | Cele-brity | Scene | Land-mark | Art-work | OCR | Total | Comm-onsense | Num-erical | Text | Code | Total |
| LLaVA-1.5 (7B) | Baseline | 190.0 | 108.3 | 96.7 | 135.0 | 140.8 | 134.7 | 153.5 | 140.0 | 75.3 | 92.5 | 1266.8 | 93.6 | 77.5 | 45.0 | 45.0 | 261.1 |
| | SPIN | 190.0 | **113.3** | **101.7** | 135.0 | **142.9** | **137.4** | 151.8 | 138.5 | **76.0** | 107.5 | **1294.0** | 93.6 | 75.0 | **60.0** | 45.0 | **273.6** |
| LLaVA-1.5 (13B) | Baseline | 185.0 | 110.0 | 100.0 | 135.0 | 146.6 | 121.8 | 161.3 | 160.3 | 87.0 | 117.5 | 1324.4 | 110.7 | 52.5 | 90.0 | 47.5 | 300.7 |
| | SPIN | 185.0 | 110.0 | 100.0 | 135.0 | 144.2 | 121.5 | **162.0** | 158.8 | **90.8** | 117.5 | **1324.7** | **112.1** | **62.5** | 87.5 | 47.5 | **309.6** |
| Qwen2.5-VL (7B) | Baseline | 180.0 | 120.0 | 148.3 | 190.0 | 169.4 | 117.9 | 139.5 | 123.3 | 141.0 | 192.5 | 1521.9 | 110.7 | 155.0 | 185.0 | 140.0 | 590.7 |
| | SPIN | **185.0** | **125.0** | 148.3 | 190.0 | **173.1** | **135.6** | 137.5 | **131.3** | 134.0 | 177.5 | **1537.3** | 110.0 | 125.0 | 185.0 | 132.5 | 552.5 |

Table 5: Results on MME evaluation. SPIN improves performance across both perception and cognition subsets on all three models. A small degradation in cognitive performance is observed for Qwen2.5-VL.

| Model | Mode | Art & Design | Business | Science | Health & Medicine | Human. & Social Sci. | Tech & Eng. | Overall |
|---|---|---|---|---|---|---|---|---|
| LLaVA-1.5 (7B) | Baseline | 51.7 | 23.3 | 24.7 | 34.0 | 49.2 | 31.4 | 34.4 |
| | SPIN | 50.0 | **24.0** | **25.3** | **35.3** | 47.5 | 31.4 | 34.4 |
| LLaVA-1.5 (13B) | Baseline | 51.7 | 23.3 | 28.0 | 39.3 | 52.5 | 32.4 | 36.6 |
| | SPIN | 51.7 | 23.3 | **29.3** | 38.0 | **54.2** | **32.9** | **36.9** |
| LLaVA-Next (7B) | Baseline | 52.5 | 30.0 | 20.7 | 36.7 | 57.5 | 26.2 | 35.3 |
| | SPIN | 52.5 | 29.3 | 20.7 | 34.0 | 57.5 | **29.5** | **35.6** |
| Qwen2.5-VL (7B) | Baseline | 68.3 | 40.0 | 40.0 | 55.3 | 67.5 | 39.5 | 49.9 |
| | SPIN | 67.5 | **42.7** | 37.3 | 54.7 | **68.3** | **40.0** | 49.9 |

Table 6: Performance on MMMU validation set for evaluating general capability. SPIN maintains or slightly improves general multimodal reasoning capabilities across all evaluated LVLMs. Following the original MMMU (Yue et al., 2024) evaluation protocol, the overall score is computed as a weighted average across all subtasks.

| Method | LLaVA-1.5 (7B) | | LLaVA-1.5 (13B) | | Qwen-VL | |
|---|---|---|---|---|---|---|
| | A | D | A | D | A | D |
| Baseline | 6.50 | 6.83 | 6.53 | 6.71 | 6.56 | 6.87 |
| SPIN | 7.19 | 6.12 | 7.03 | 6.18 | 7.08 | 6.22 |

Table 7: GPT-4o assisted evaluation showcasing Accuracy (A) and Detailedness (D) of the generated captions.

(D): measuring the richness of necessary details in the responses. The prompt we used is given in the Appendix Table 14. It is designed to ensure that the sequential order in which the responses are presented does not affect the judgement. The evaluation model is prompted to identify objects misaligned with the image context, identifying any discrepancies in count, position, and colors of objects in the images. The results presented in Table 7 show that SPIN improves Accuracy up to ∼7%, while maintaining Detailedness close to baseline. We use the same values for Layers, $r$, and $\alpha$ reported for CHAIR for all captioning tasks.

## 4.6 General Performance Assessment

**MME Evaluation:** In Table 5, we present SPIN results on MME full set (Fu et al., 2023) that assesses the perceptual and cognitive abilities of VLMs across a total of 14 subtasks, including tasks such as OCR, visual knowledge, attribute relationships, and object recognition. SPIN outperforms the base-line on both perception and cognition tasks for all three models, with a small degradation in cognition for Qwen2.5-VL. The improvements in the first four subtasks relating to existence, count, position, and colour demonstrates the model's ability to reduce object-level and attribute-level hallucinations.

**MMMU Evaluation:** In Table 6, we present MMMU (Yue et al., 2024) (a standard multimodal understanding and reasoning benchmark) scores for SPIN to assess its impact on the general capability of LVLMs. SPIN attains the same performance as baseline for LLaVA-1.5 (7B) and Qwen2.5-VL, while attaining a small improvement on LLaVA-1.5 (13B) and LLaVA-Next. This demonstrates that head suppression does not hinder, but enhances the model's core capabilities while providing substantial hallucination reduction. Notably, we use the same parameters for head suppression (Layers, $r$, and $\alpha$) as reported for CHAIR evaluation for both MME and MMMU.

## 4.7 Throughput Estimation

To assess the efficiency of different hallucination mitigation strategies, we compare their throughput in Figure 6. Throughput is computed by measuring the total computational latency for generating tokens, and dividing the total number of generated

### LLaVA-1.5 (7B)

| Layers | $r$ | $\mathbf{C}_S$ (↓) | $\mathbf{C}_I$ (↓) | F1 (↑) |
|---|---|---|---|---|
| 1~32 | 0.05 | **17.4** | **5.8** | 74.6 |
| 1~32 | 0.10 | 22.8 | 7.4 | 74.1 |
| 1~32 | 0.15 | 20.8 | 6.4 | 73.8 |
| 1~16 | 0.05 | 18.2 | 6.2 | 73.7 |
| 1~20 | 0.05 | 18.4 | 6.3 | 74.1 |
| 1~24 | 0.05 | 19.2 | 6.5 | 73.8 |
| 4~32 | 0.05 | 42.0 | 11.4 | 78.2 |
| 8~32 | 0.05 | 42.4 | 11.1 | **78.8** |

### LLaVA-1.5 (13B)

| Layers | $r$ | $\mathbf{C}_S$ (↓) | $\mathbf{C}_I$ (↓) | F1 (↑) |
|---|---|---|---|---|
| 1~32 | 0.05 | 45.8 | 12.0 | 78.2 |
| 1~32 | 0.10 | 29.6 | 8.0 | 78.9 |
| 1~32 | 0.15 | 41.6 | 10.3 | **80.2** |
| 1~32 | 0.20 | 36.0 | 10.0 | 79.1 |
| 1~16 | 0.10 | 30.6 | 8.3 | 79.6 |
| 1~20 | 0.10 | 29.2 | **7.9** | 79.1 |
| 1~24 | 0.10 | **28.0** | 9.2 | 79.3 |

### Qwen-VL

| Layers | $r$ | $\mathbf{C}_S$ (↓) | $\mathbf{C}_I$ (↓) | F1 (↑) |
|---|---|---|---|---|
| 1~32 | 0.05 | 44.4 | 11.7 | **77.6** |
| 1~32 | 0.10 | 39.8 | 10.9 | 77.1 |
| 1~32 | 0.15 | 41.6 | 11.7 | 77.5 |
| 1~32 | 0.20 | 39.6 | 11.2 | 76.9 |
| 1~32 | 0.25 | 36.4 | 9.7 | 77.3 |
| 1~32 | 0.30 | **26.0** | 9.1 | 75.6 |
| 1~20 | 0.30 | 27.6 | **7.9** | 76.1 |
| 1~24 | 0.30 | 28.4 | 8.3 | 76.2 |
| 8~32 | 0.30 | 32.4 | 9.4 | 76.3 |

Table 8: Ablation on the ratio of suppressed heads and layer selection for three different models. The best $r$ is identified in the top part of the tables (marked with blue ). The layer configuration for selected value of $r$ is presented in the bottom part (best layer configuration is marked in yellow ).
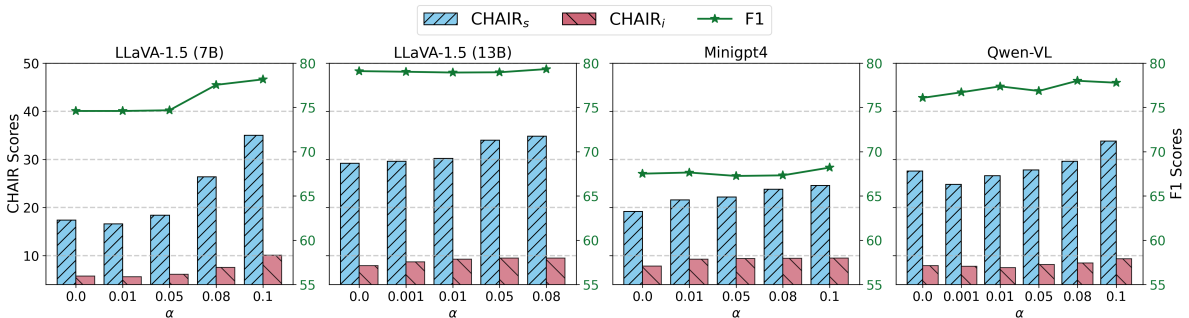


Figure 5: Ablation on the suppression factor ($\alpha$) for four models.

tokens by this latency. SPIN achieves the highest throughput among SoTA approaches, while achieving throughput performance close to baseline.
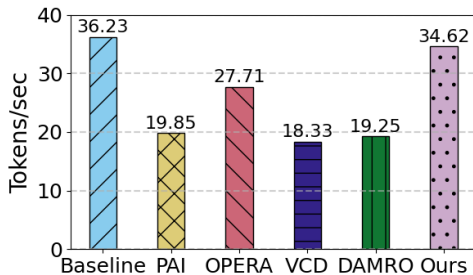


Figure 6: Throughput comparison with existing methods given by the number of tokens generated per second.

### 4.8 Ablation Studies

We present the hyperparameter evaluations for SPIN in Table 8, where we first vary $r$ to obtain reduced hallucinations without a significant drop in F1. We obtain $r = 0.05$ for LLaVA 7B, 0.10 for LLaVA 13B, and 0.30 for Qwen-VL. Next, we evaluate different layer configurations to find hallucination-prone layers . While uniformly pruning across all layers works best for LLaVA 7B, initial layers are more hallucination-prone for LLaVA

13B and Qwen-VL. Finally, for selected $r$ and layers for head suppression, we plot $C_s$, $C_i$ and F1 for different $\alpha$ in Figure 5. F1 scores show comparable values for different $\alpha$ for MiniGPT-4 and LLaVA 13B models. However, for LLaVA 7B and Qwen-VL, complete pruning leads to a drop in F1. This happens because the pre-trained model weights are not adapted to this change. Therefore, we perform partial head suppression using the parameter $\alpha$, where a higher value of $\alpha$ preserves the F1 score by reducing the suppression intensity.

## 5 Conclusions

Despite achieving remarkable progress in vision-language tasks, LVLMs still suffer from hallucinations arising from misalignment between visual input and the generated text. We identify a dynamic set of attention heads based on input query as the potential cause for hallucinations. To counteract this, we propose an image-attention guided head pruning strategy agnostic to decoding method, that is run-time efficient and can be seamlessly integrated during inference. Our approach leads to a substantial reduction in hallucinations, outperforming existing methods while maintaining accuracy.

## Limitations

We propose an inference-only solution for ease of integration and to accommodate scenarios with limited training resources. However, our approach could be enhanced by incorporating a trainable router for more adaptive head selection. Our method's effectiveness is somewhat reduced when paired with sampling-based decoding strategies (e.g., nucleus sampling). We attribute this to the greater token selection variability inherent in these stochastic approaches (as further detailed in Section 4.3). Our method requires access to the model weights, which limits its applicability to open-source models. It cannot be applied to closed-source or API-based models, where direct access to weights is restricted.

**Usage and License:** We acknowledge the licensing terms associated with the artifacts used in this study. LLaVA (Liu et al., 2024b) is released under an open-source license (Apache 2.0), permitting modification and distribution with attribution. MiniGPT-4 (Zhu et al., 2024) follows similar licensing terms under the LLaMA model's restrictions, where the base model weights are subject to Meta's research agreement. Qwen-VL models are also released under Apache 2.0, whereas Shikra is released under the Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0) license. Our use of these models and datasets is strictly for academic research purposes, adhering to all licensing and usage guidelines. All datasets used in our experiments comply with their respective licenses, and we ensure proper attribution where required.

## References

2024. Shiftaddllm: Accelerating pretrained llms via post-training multiplication-less reparameterization. *Advances in Neural Information Processing Systems*, 37.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. 2025. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*.

Junzhe Chen, Tianshu Zhang, Shiyu Huang, Yuwei Niu, Linfeng Zhang, Lijie Wen, and Xuming Hu. 2025. Ict: Image-object cross-level trusted intervention for mitigating object hallucination in large vision-language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, pages 4209–4221.

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

Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. 2024a. An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models. In *European Conference on Computer Vision*, pages 19–35. Springer.

Zhaorun Chen, Zhuokai Zhao, Hongyin Luo, Huaxiu Yao, Bo Li, and Jiawei Zhou. 2024b. Halc: Object hallucination reduction via adaptive focal-contrast decoding. *arXiv preprint arXiv:2403.00425*.

Markus Freitag and Yaser Al-Onaizan. 2017. Beam search strategies for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 56–60, Vancouver. Association for Computational Linguistics.

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. 2023. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*.

Xuan Gong, Tianshi Ming, Xinpeng Wang, and Zhihua Wei. 2024. 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. Association for Computational Linguistics.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *International Conference on Learning Representations*.

Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. 2024. Opera: Alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13418–13427.

Jitesh Jain, Jianwei Yang, and Humphrey Shi. 2024. Vcoder: Versatile vision encoders for multimodal

large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27992–28002.

Peng Jin, Bo Zhu, Li Yuan, and Shuicheng Yan. 2024. Moh: Multi-head attention as mixture-of-head attention. *arXiv preprint arXiv:2410.11842*.

Liqiang Jing and Xinya Du. 2025. Fgaif: Aligning large vision-language models with fine-grained ai feedback. *Preprint*, arXiv:2404.05046.

Hao Kang, Qingru Zhang, Souvik Kundu, Geonhwa Jeong, Zaoxing Liu, Tushar Krishna, and Tuo Zhao. 2024. Gear: An efficient kv cache compression recipefor near-lossless generative inference of llm. *NuerIPS Workshop ESNLP*.

Ivan Krasin, Tom Duerig, Neil Alldrin, Andreas Veit, Sami Abu-El-Haija, Serge Belongie, David Cai, Zheyun Feng, Vittorio Ferrari, Victor Gomes, Abhinav Gupta, Dhyanesh Narayanan, Chen Sun, Gal Chechik, and Kevin Murphy. 2016. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from https://github.com/openimages*.

Souvik Kundu, Anahita Bhiwandiwalla, Sungduk Yu, Phillip Howard, Tiep Le, Sharath Nittur Sridhar, David Cobbley, Hao Kang, and Vasudev Lal. 2025. Lvlm-compress-bench: Benchmarking the broader impact of large vision-language model compression. *Annual Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics*.

Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. 2024. 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*, pages 13872–13882.

Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer.

Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2023. Mitigating hallucination in large multi-modal models via robust instruction tuning. In *The Twelfth International Conference on Learning Representations*.

Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li, and Wei Peng. 2024a. A survey on hallucination

in large vision-language models. *arXiv preprint arXiv:2402.00253*.

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024b. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 26296–26306.

Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024c. Llava-next: Improved reasoning, ocr, and world knowledge.

Sheng Liu, Haotian Ye, and James Zou. 2025. Reducing hallucinations in large vision-language models via latent space steering. In *The Thirteenth International Conference on Learning Representations*.

Shi Liu, Kecheng Zheng, and Wei Chen. 2024d. Paying more attention to image: A training-free method for alleviating hallucination in lvlms. In *European Conference on Computer Vision*, pages 125–140. Springer.

Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2018. Object hallucination in image captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4035–4045, Brussels, Belgium. Association for Computational Linguistics.

Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liangyan Gui, Yu-Xiong Wang, Yiming Yang, Kurt Keutzer, and Trevor Darrell. 2024. Aligning large multimodal models with factually augmented RLHF. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 13088–13110, Bangkok, Thailand. Association for Computational Linguistics.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. 2019. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. *arXiv preprint arXiv:1905.09418*.

Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *CVPR*.

Haozhe Zhao, Shuzheng Si, Liang Chen, Yichi Zhang, Maosong Sun, Mingjia Zhang, and Baobao Chang. 2024. Looking beyond text: Reducing language bias

in large vision-language models via multimodal dual-attention and soft-image guidance. *arXiv preprint arXiv:2411.14279*.

Wenhao Zheng, Yixiao Chen, Weitong Zhang, Souvik Kundu, Yun Li, Zhengzhong Liu, Eric P Xing, Hongyi Wang, and Huaxiu Yao. 2025. Citer: Collaborative inference for efficient large language model decoding with token-level routing. *arXiv preprint arXiv:2502.01976*.

Zhenhong Zhou, Haiyang Yu, Xinghua Zhang, Rongwu Xu, Fei Huang, Kun Wang, Yang Liu, Junfeng Fang, and Yongbin Li. 2024. On the role of attention heads in large language model safety. *arXiv preprint arXiv:2410.13708*.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2024. MiniGPT-4: Enhancing vision-language understanding with advanced large language models. In *The Twelfth International Conference on Learning Representations*.

# A  Appendix

## A.1  Baseline Hyperparameters

The hyperparameters $\alpha$ and $r$ for SPIN are specified in each table. For existing methods, we use the best-performing configurations reported in their respective papers.

For PAI, we adopt the following hyperparameter settings: for LLaVA, we use $\alpha = 0.5$, $\gamma = 1.2$, and $layer\ prior = 2$; for MiniGPT-4, $\alpha = 0.2$, $\gamma = 1.3$, and $layer\ prior = 3$; and for Shikra, $\alpha = 0.6$, $\gamma = 1.1$, and $layer\ prior = 3$. Since the original PAI paper does not report results on Qwen-VL, we experimented with several configurations and found that reusing Shikra's settings yielded the best performance.

For OPERA, we use the same configuration across all models: $\sigma = 50$, $r = 15$, $N_{\text{can}} = 5$, and $\alpha = 1$. For VCD, we use $\alpha = 1$, $\beta = 0.1$, and apply noise steps $T = 999$ for all models. For DAMRO, we select the top 10 tokens as outliers and set $\alpha = 0.5$.

## A.2  Ablation Studies

**Head Selection Strategy:** We explore a variety of training-free head selection strategies, as summarized in Table 9. Query-based head selection proposed in (Jin et al., 2024) is computed based on the L2 norm of the query tokens, whereas the key based selection is performed based on the L2 norm of the key tokens. In addition, we experiment with head selection guided by the overall attention (to both text and image tokens) as well as attention to the image tokens only. Among these strategies, selecting heads based solely on image attention proves to be the most effective for hallucination reduction, yielding the best CHAIR scores. Here, we use $\alpha = 0$, completely pruning attention heads across all layers.

| Selection Strategy | Supp. Heads | CHAIR$_S$ ($\downarrow$) | CHAIR$_I$ ($\downarrow$) | F1 ($\uparrow$) |
|---|---|---|---|---|
| Query | 0.10 | 40.4 | 11.6 | 77.4 |
| Query | 0.25 | 34.6 | 10.4 | 75.8 |
| Key | 0.25 | 44.6 | 12.4 | 77.1 |
| Attention | 0.15 | 20.4 | 8.0 | 73.4 |
| Attention | 0.05 | 17.6 | 6.1 | 74.5 |
| Image Attention | 0.05 | **17.4** | **5.8** | 74.6 |

Table 9: Comparison of different training-free head selection strategies in SPIN on LLaVA-1.5 (7B) using CHAIR evaluation. Selecting heads based on image attention yields the most effective reduction in hallucination.

| Model | Mode | Layers | $r$ | $\alpha$ | Accuracy | F1 |
|---|---|---|---|---|---|---|
| LLaVA-1.5 (7B) | Baseline | - | - | - | 85.69 | 84.28 |
| | SPIN | 1~16 | 0.15 | 0.0 | 87.84 | 87.13 |
| | SPIN | 1~16 | 0.15 | 0.001 | 87.79 | 87.07 |
| | SPIN | 1~32 | 0.20 | 0.08 | **88.19** | **87.42** |
| LLaVA-1.5 (13B) | Baseline | - | - | - | 85.12 | 84.48 |
| | SPIN | 1~32 | 0.35 | 0.0 | 88.47 | 88.00 |
| | SPIN | 1~24 | 0.35 | 0.0 | 88.40 | 88.01 |
| | SPIN | 1~20 | 0.35 | 0.001 | **88.50** | **88.17** |
| MiniGPT-4 | Baseline | - | - | - | 78.08 | 76.20 |
| | SPIN | 1~32 | 0.05 | 0.0 | **79.43** | **78.78** |
| | SPIN | 1~24 | 0.05 | 0.08 | 79.37 | 78.61 |
| | SPIN | 1~24 | 0.05 | 0.001 | 79.14 | 78.46 |
| Shikra | Baseline | - | - | - | 78.27 | 78.89 |
| | SPIN | 1~16 | 0.20 | 0.0 | **79.20** | 79.40 |
| | SPIN | 1~16 | 0.20 | 0.01 | 79.17 | 79.37 |
| | SPIN | 1~16 | 0.20 | 0.05 | 79.19 | **79.45** |

Table 10: Multi-turn POPE evaluation under greedy decoding. SPIN consistently reduces hallucinations across all models on multi-turn vision-language tasks.

## A.3  Head Mask Visualization

The head masks for POPE and CHAIR using 95% active heads are visualized in Figure 7. The pruned heads are more consistent for POPE, likely due to the fact that "Yes" or "No" answers are generated as the first token. In contrast, CHAIR generates longer sequences, causing tokens at each position to likely prune different heads.



Figure 7: Heatmap for head masks averaged over all test samples for CHAIR (left) and POPE random split (right) for 95% active heads. White and dark blue indicate the always pruned and always active heads, respectively.

## A.4  Extended Evaluation of SPIN

**POPE Evaluation:** We assess SPIN on POPE evaluation to examine its performance in multi-turn dialogue settings. All experiments in this section are performed using greedy decoding. As shown in Table 10, SPIN consistently reduces hallucinations across multiple models. These results demonstrate that our method not only generalizes well beyond CHAIR, but also achieves strong hallucination mitigation on multi-turn vision-language tasks.

**CHAIR Results under Diverse Decoding Strategies:** To further evaluate the effectiveness of SPIN, we conduct experiments on CHAIR under various decoding strategies: greedy search, beam

| Model | Decoding | Method | Layers | Supp. Heads | Scale Factor | CHAIR$_S$ (↓) | CHAIR$_I$ (↓) | F1 (↑) |
|---|---|---|---|---|---|---|---|---|
| MiniGPT-4 | Beam Search | Baseline | - | - | - | 33.4 | 11.0 | **70.2** |
| | | OPERA | - | - | - | 23.2 | 8.8 | 69.4 |
| | | PAI | - | - | - | 22.2 | 8.3 | 68.1 |
| | | SPIN | 1∼24 | 0.15 | 0.0 | 16.6 | **4.8** | 68.2 |
| | | SPIN | 1∼16 | 0.18 | 0.05 | **16.2** | 6.6 | 69.1 |
| | Nucleus Sampling | Baseline | - | - | - | 31.8 | 11.8 | 66.4 |
| | | PAI | - | - | - | **23.0** | **9.0** | **68.3** |
| | | SPIN | 1∼16 | 0.18 | 0.05 | 23.6 | 12.5 | 62.9 |
| Shikra | Beam Search | Baseline | - | - | - | 52.6 | 14.1 | **75.7** |
| | | OPERA | - | - | - | 40.2 | 12.6 | 72.7 |
| | | PAI | - | - | - | 38.6 | 10.1 | 74.9 |
| | | SPIN | 1∼32 | 0.40 | 0.0 | **32.4** | **10.0** | 74.5 |
| | Nucleus Sampling | Baseline | - | - | - | 56.0 | 15.7 | 72.8 |
| | | VCD | - | - | - | 55.6 | 15.2 | 74.9 |
| | | PAI | - | - | - | **33.8** | **8.7** | **75.1** |
| | | SPIN | 1∼32 | 0.40 | 0.0 | 34.2 | 10.6 | 71.5 |

Table 11: CHAIR evaluation using MiniGPT-4 and Shikra under beam search and nucleus sampling decoding.

| Model | Method | Layers | $r$ | $\alpha$ | C$_S$(↓) | C$_I$(↓) | F1(↑) |
|---|---|---|---|---|---|---|---|
| LLaVA-Next (7B) | Baseline | - | - | - | 35.6 | 8.6 | 71.7 |
| | SPIN | 1∼32 | 0.05 | 0.0 | 32.6 | 8.1 | **72.3** |
| | SPIN | 1∼16 | 0.20 | 0.05 | **26.8** | **7.4** | 71.7 |
| Qwen2.5-VL (7B) | Baseline | - | - | - | 34.8 | 7.9 | **75.7** |
| | SPIN | 1∼16 | 0.02 | 0.0 | **28.6** | **7.0** | 74.0 |
| | SPIN | 1∼16 | 0.02 | 0.05 | 30.6 | 7.5 | 74.7 |

Table 12: CHAIR evaluation results on the latest LVLMs, including LLaVA-Next and Qwen2.5-VL.

| Model | Mode | Layers | $r$ | $\alpha$ | Accuracy | F1 |
|---|---|---|---|---|---|---|
| LLaVA-Next (7B) | Baseline | - | - | - | 88.45 | 87.99 |
| | SPIN | 1∼32 | 0.05 | 0.0 | 89.00 | 88.52 |
| | SPIN | 1∼32 | 0.05 | 0.001 | **89.07** | **88.60** |
| Qwen2.5-VL (7B) | Baseline | - | - | - | 81.64 | 77.75 |
| | SPIN | 1∼20 | 0.02 | 0.0 | 83.02 | 79.80 |
| | SPIN | 1∼20 | 0.02 | 0.001 | **83.05** | **79.87** |

Table 13: Experimental results on Multi-turn POPE using LLaVA-Next and Qwen2.5-VL.

search, and nucleus sampling. Across four models: LLaVA-1.5 (7B), MiniGPT-4, Shikra, and Qwen-VL. We explore a wider range of configurations by varying the ratio of suppressed heads, scaling factors, and the layers where SPIN is applied.

As discussed in Section 4.2, we have already provided a detailed analysis of SPIN's performance under greedy search for all four models, as well as its behavior on beam search and nucleus sampling for LLaVA-1.5 (7B) and Qwen-VL. Here, we highlight the results for MiniGPT-4 and Shikra under beam search and nucleus sampling, as shown in Table 11. For MiniGPT-4, SPIN reduces CHAIR$_I$ and CHAIR$_S$ by 2.3× and 2.0× with beam search. For Shikra, the corresponding reductions are 1.4× and 1.6×. Under nucleus sampling, SPIN yields less significant improvements, but still effectively mitigates hallucinations. For MiniGPT-4, SPIN reduces CHAIR$_S$ by 1.35×. On Shikra, it reduces CHAIR$_I$ and CHAIR$_S$ by 1.5× and 1.6×, with only a 1.3% drop in F1.

## A.5 Evaluations on the Latest Models

To further assess the accessibility and robustness of our method on state-of-the-art LVLMs, we con-duct evaluations under greedy decoding on LLaVA-Next (Liu et al., 2024c) and Qwen2.5-VL (Bai et al., 2025). The CHAIR results are reported in Table 12, while the POPE results are shown in Table 13. SPIN consistently alleviates hallucinations in these newer models and provides improved accuracy and F1 for POPE evaluation.

It is worth noting that, due to differences in architectures in the language backbone, Qwen2.5-VL generally benefits more from suppressing a smaller fraction of heads. Specifically, Qwen2.5-VL comprises 28 layers with 28 heads per layer, whereas LLaVA-Next, leveraging Mistral-7B, consists of 32 layers with 32 heads each. This indicates that each attention head in Qwen2.5-VL undertakes greater functional responsibility, so overly aggressive suppression of heads can undermine the model's capabilities.

## A.6 CHAIR Visualization

The captions generated using the baseline and SPIN are visualized in Figures 8, 9, 10, 11, and 12, for LLaVA-1.5 (7B), LLaVA-1.5 (13B), Shikra, Qwen-VL, and MiniGPT-4. SPIN effectively mitigates hallucinations generated by the baseline models.

Table 14: The prompt used for GPT-4o evaluation adopted from Liu et al. (2024d).

**GPT-4o Prompt**

You are required to score the performance of two AI assistants in describing a given image. You should pay extra attention to the hallucination, which refers to the part of descriptions that are inconsistent with the image content, such as claiming the existence of something not present in the image or describing incorrectly in terms of the counts, positions, or colors of objects in the image. Please rate the responses of the assistants on a scale of 1 to 10, where a higher score indicates better performance, according to the following criteria:

1: Accuracy: whether the response is accurate with respect to the image content. Responses with fewer hallucinations should be given higher scores.

2: Detailedness: whether the response is rich in necessary details. Note that hallucinated descriptions should not count as necessary details.

Please output the scores for each criterion, containing only two values indicating the scores for Assistant 1 and 2, respectively. The two scores are separated by a space. Following the scores, please provide an explanation of your evaluation, avoiding any potential bias and ensuring that the order in which the responses were presented does not affect your judgment.

[Assistant 1]
{Response of Assistant 1}
[End of Assistant 1]

[Assistant 2]
{Response of Assistant 2}
[End of Assistant 2]

Output format:
Accuracy: <Scores of the two answers>
Reason:

Detailedness: <Scores of the two answers>
Reason:

**Instruction:** Please help me describe the image in detail.

**LLaVA-1.5-7B:**

The image features a man playing a guitar in a bedroom, standing on a bed and playing with great enthusiasm. He is holding the guitar with both hands, fully engaged in his performance. The bed is the main focus of the scene, with the man standing on it while playing.

In the room, there are several books scattered around, some on the bed and others on the floor. A bottle can also be seen on the floor, possibly containing a beverage for the man to enjoy during his guitar playing session.

**SPIN:**

The image features a man playing a guitar in a bedroom, standing on a bed and playing the guitar with great enthusiasm. The man is wearing a black shirt and appears to be enjoying the moment.

The room has a bed with a purple blanket, and a few books scattered around the area. There is also a bottle placed near the bed, possibly containing a beverage. The overall atmosphere of the room suggests a casual and relaxed environment.

**Input Image**

---

**Instruction:** Please help me describe the image in detail.

**LLaVA-1.5-7B:**

The image features a dining table with a red placemat underneath it. On the table, there are several jars filled with oranges, some of which are whole and others are sliced. The oranges are arranged in various positions, with some jars placed closer to the edge of the table and others in the middle. The table also has a bowl and a cup, adding to the overall presentation of the scene.

**SPIN:**

The image features a dining table with a red placemat underneath it. On the table, there are several jars filled with sliced oranges, creating a visually appealing display. The oranges are arranged in various positions, with some jars placed closer to the front and others further back. The scene gives off a fresh and vibrant atmosphere, showcasing the oranges in a visually pleasing manner.

**Input Image**

---

**Instruction:** Please help me describe the image in detail.

**LLaVA-1.5-7B:**

The image features a wooden closet with numerous ties hanging on the wall. There are at least 13 ties of various colors and patterns, including blue, green, and red. The ties are arranged in a neat and orderly manner, creating an organized display.

In addition to the ties, there are two handbags placed on the floor, one on the left side and the other on the right side of the closet. A suitcase can also be seen in the background, possibly indicating a recent trip or a planned journey.

**SPIN:**

The image features a well-organized collection of ties hanging neatly on a wall. There are numerous ties of various colors and patterns, creating a visually appealing display. The ties are arranged in a row, with some ties hanging higher than others, giving the impression of a well-curated collection. The ties are displayed in a room, possibly a closet or a store, showcasing the variety of options available for customers or for personal use.

**Input Image**

Figure 8: SPIN's effectiveness in reducing hallucinations with LLaVA-1.5 (7B).

**Instruction:** Please help me describe the image in detail.

**LLaVA-1.5-13B:**

The image features two dogs lying on a bed, one on the left side and the other on the right side. They appear to be relaxed and comfortable, possibly sleeping or resting. The bed is covered with a blue blanket, providing a cozy environment for the dogs.

In addition to the dogs, there is a chair located in the background, positioned behind the bed. The scene captures a peaceful and restful moment shared between the two dogs.

**SPIN:**

The image features two dogs lying on a bed, with one dog being a white dog and the other being a black dog. The white dog is lying on the left side of the bed, while the black dog is lying on the right side. Both dogs appear to be sleeping or resting, possibly enjoying each other's company. The bed is covered with a blanket, providing a comfortable surface for the dogs to lay on.

**Input Image**

---

**Instruction:** Please help me describe the image in detail.

**LLaVA-1.5-13B:**

The image features a man standing in a body of water, interacting with a large elephant. The man is petting the elephant's trunk, and the elephant appears to be enjoying the attention. The scene takes place in a natural environment, with the man and the elephant surrounded by water.

There are several other people in the scene, some of them standing closer to the water's edge, while others are further away. The presence of multiple people suggests that this might be a social gathering or an event involving the elephant.

**SPIN:**

The image features a man and an elephant in a body of water. The man is standing in the water, washing the elephant, while the elephant is standing in the water as well. The man is wearing a white shirt and appears to be enjoying the interaction with the elephant. The scene is set in a natural environment, with the elephant and the man surrounded by water.

**Input Image**

---

**Input Image**

**Instruction:** Please help me describe the image in detail.

**LLaVA-1.5-13B:**

The image features two men standing under a red tent, exchanging a medal. One man is wearing a suit and tie, while the other is dressed in a more casual outfit. They are both holding the medal, with one man handing it to the other.

The scene takes place outdoors, with a dining table visible in the background. There are also two chairs placed near the table, and a potted plant can be seen in the foreground. The men appear to be engaged in a formal or ceremonial event, as they are exchanging the medal under the tent.

**SPIN:**

The image features two men standing under a red tent, exchanging a medal. One man is wearing a suit and tie, while the other man is wearing a white shirt. They are both holding a medal, and it appears that one man is giving the medal to the other. The men are standing close to each other, with one man on the left and the other on the right.

Figure 9: SPIN's effectiveness in reducing hallucinations with LLaVA-1.5 (13B).

**Instruction:** Please help me describe the image in detail.

**Shikra:**

The image features a small, adorable pug dog sitting on a skateboard in a room. The dog is wearing a collar, and it appears to be wearing a harness. The skateboard is placed on a wooden floor, and the pug is positioned in the center of the board, looking rather uncomfortable or perhaps curious.

There are a few additional details in the room, such as a refrigerator on the right side of the scene, a chair near the left wall, and another chair in the middle of the room.

**SPIN:**

The image features a small, adorable pug sitting on a skateboard, likely on a wooden floor. The dog appears to be wearing a collar, indicating that it is a pet. The scene is set in a room with a yellow wall, giving it a warm and inviting atmosphere. The dog is sitting on the skateboard, showcasing a sense of confidence and authority in its positioning.

**Input Image**

---

**Instruction:** Please help me describe the image in detail.

**Shikra:**

The image features a street scene with a prominent orange detour sign placed on the road. The sign is large and clearly displays the direction to the detour, likely due to a road closure or a detour route.

There are several cars in the scene, with one car on the left side of the image, another near the center, and two more cars on the right side of the picture. Additionally, a traffic light can be seen in the background, further indicating that this is a busy area with traffic regulations.

**SPIN:**

The image features a street scene with a prominent orange detour sign placed on the side of the road. The sign is large and clearly visible, giving the impression that the road is a major detour. The sign is pointing to the left and indicates that the road is closed or unavailable.

The image also shows that the road is covered with snow, which suggests that the snowfall has recently occurred. The presence of snow on the road indicates that the weather is cold and snowy, and that the road is closed due to winter conditions.

**Input Image**

---

**Instruction:** Please help me describe the image in detail.

**Shikra:**

The image features a small piece of cake served on a plate, accompanied by a cup of milk. The cake has a bite taken out of it, and the milk is partially spilled, creating a mess. A fork is placed next to the cake, suggesting that someone is about to enjoy the treat. The scene is set on a dining table, making for a simple and casual snack.

**SPIN:**

The image features a small piece of cake served on a plate, placed next to a fork and a cup of milk. The cake has a bite taken out of it, indicating that someone has started eating it. The plate and fork are positioned on a dining table, creating an inviting scene. The presence of the fork signifies that the person consuming the meal is ready to be served.

**Input Image**

Figure 10: SPIN's effectiveness in reducing hallucinations with Shikra.

**Instruction:** Describe this image in detail.

**Qwen-VL-Chat:**

The image features a serene body of water with a pair of geese floating peacefully on it. Two large boats are docked on the water, one located closer to the left side and the other on the right side. The boats appear to be cruise ships or similar in nature, as they are quite large compared to the small boats often seen on waterways.

The scene is set against a beautiful sunset, with the clouds and sky providing a stunning backdrop for the boats and the water. The geese add a touch of nature and tranquility to the scene, making it a picturesque representation of a leisurely waterfront.

**SPIN:**

The image features a serene body of water with a pair of large boats floating peacefully along the water's surface. In total, there are four ducks swimming around the boats, enjoying their time on the water. The ducks are scattered, with some swimming closer to the boats and others floating further away, exploring the lake.

The sky above is filled with a variety of clouds, adding a dynamic element to the scene. The combination of the boats and ducks creates a tranquil atmosphere, as if everyone is enjoying a leisurely day on the water.

---

**Instruction:** Describe this image in detail.

**Qwen-VL-Chat:**

The image features a large boat floating on the ocean, with the word \"EUKOR\" written on it. There is another smaller boat nearby, which appears to be a barge. In the sky, a couple of airplanes can be seen, one high up and another one lower down.

There are two people visible on the boats, likely the crew or passengers. The scene captures the essence of a large marine operation or a search and rescue mission in progress.

**SPIN:**

The image features a large boat floating on the ocean, with a sky background filled with airplanes. There are two airplanes flying at different altitudes over the ocean, with one plane flying higher in the sky and another one lower in the sky. Additionally, there is a helicopter flying relatively closer to the boat. The boat is named Eukor, and it is positioned in the lower part of the image, with the ocean stretching out behind it.



**Input Image**

---



**Input Image**

**Instruction:** Describe this image in detail.

**Qwen-VL-Chat:**

The image features a tall skyscraper building all lit up at night, standing tall in the background. In the foreground, there is a traffic light with a \"Green Light\" sign below it. The traffic light is situated on the corner of a street, with the illuminated skyscraper serving as a beautiful backdrop.

**SPIN:**

The image features a tall skyscraper building lit up with red lights at night, with a traffic light situated next to the building. The tall building has many windows and is illuminated by additional lights on the side of the building. The scene captures a street sign and traffic light, emphasizing the urban environment.
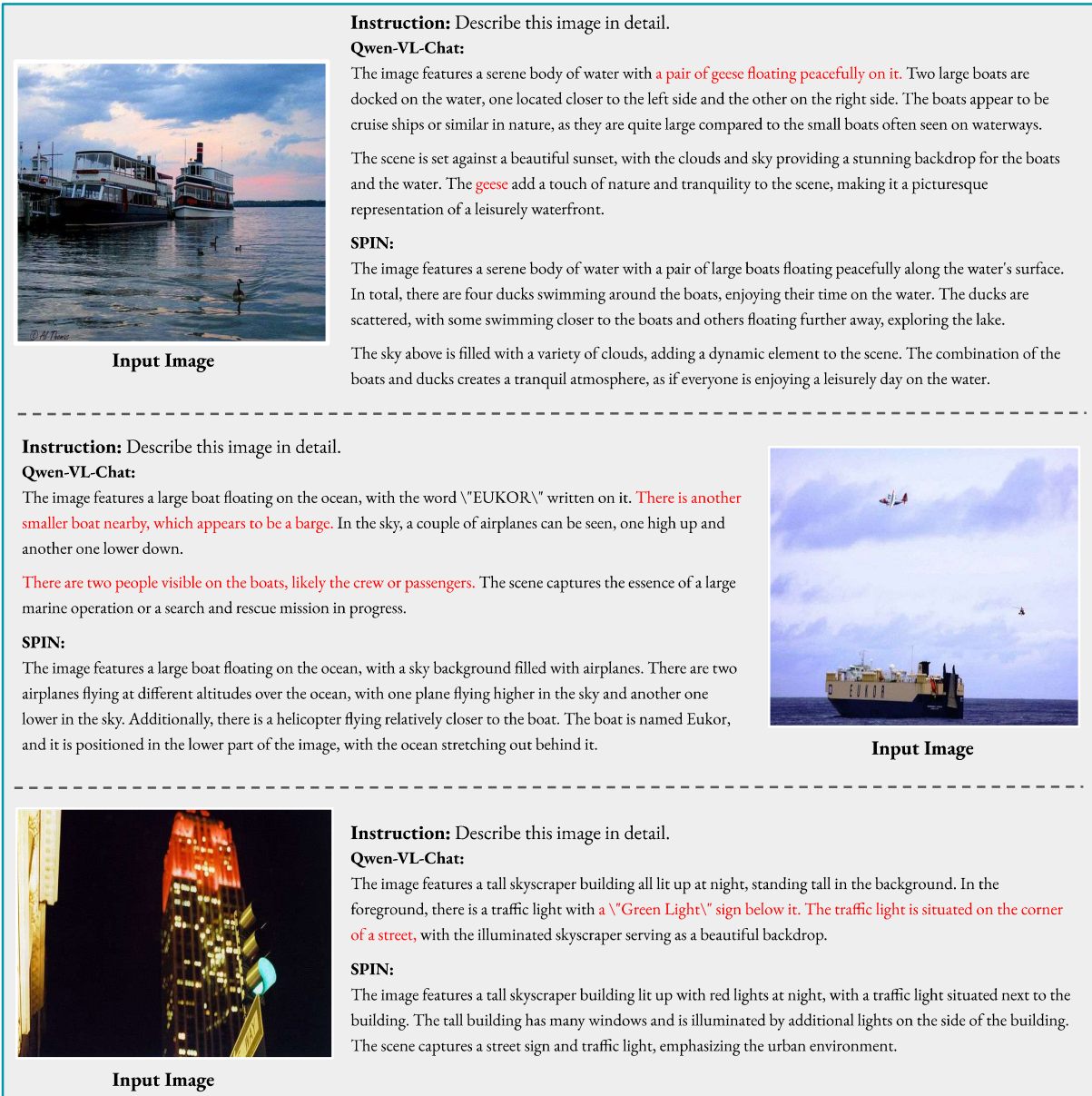
Figure 11: SPIN's effectiveness in reducing hallucinations with Qwen-VL-Chat.

**Instruction:** Please help me describe the image in detail.

**Minigpt4:**
The image shows a man standing in a green field with a brown cow next to him. The man is holding a rope and the cow is looking up at him. There are trees in the background and a blue sky with clouds.

**SPIN:**
The image shows a man standing in a green field with two cows grazing on the grass. The man is holding a rope and leading one of the cows, while the other cow is walking behind him. The sky is clear blue with some clouds visible in the distance.
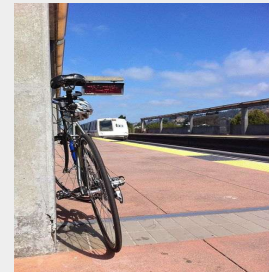
---

**Instruction:** Please help me describe the image in detail.

**Minigpt4:**
The image shows a bicycle leaning against a concrete pillar in a train station. The bicycle has a black frame and white tires. The station has a metal roof and a sign with the name of the station. There is a blue sky and clouds in the background.

**SPIN:**
The image shows a bicycle leaning against a concrete pillar in front of a train station. There is a blue sky and clouds visible in the background.

---

**Instruction:** Please help me describe the image in detail.

**Minigpt4:**
The image shows a man riding a surfboard in the ocean. He is wearing a black and red wetsuit and has a black and red helmet on his head. He is holding onto the surfboard with his hands and has a white surfboard under his arm. The ocean is turquoise and there are white waves in the background. The man is smiling and looks happy.

**SPIN:**
The image shows a man riding a wave on his surfboard in the ocean. He is wearing black wetsuit and has a red board under him. The water is green and blue with white waves crashing on the shore.
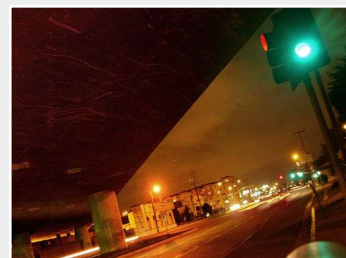
---

**Instruction:** Please help me describe the image in detail.

**Minigpt4:**
The image shows a street at night with a green traffic light on the left side of the image. The light is illuminating the road and the buildings on either side of it. The buildings are tall and have many windows, and there are streetlights on either side of the road. The sky is dark and there are stars visible in the background. The image is taken from a high angle and the view is looking down the street.

**SPIN:**
The image shows a street scene at night with green traffic light on the road and buildings in the background. The buildings are tall and have many windows, some of which are lit up. There is also a sidewalk on the left side of the road.

Figure 12: SPIN's effectiveness in reducing hallucinations with MiniGPT-4.