

# ONLY: One-Layer Intervention Sufficiently Mitigates Hallucinations in Large Vision-Language Models

Zifu Wan\* Ce Zhang\* Silong Yong Martin Q. Ma Simon Stepputtis

Louis-Philippe Morency Deva Ramanan Katia Sycara Yaqi Xie

Carnegie Mellon University

## Abstract

Recent Large Vision-Language Models (LVLMs) have introduced a new paradigm for understanding and reasoning about image input through textual responses. Although they have achieved remarkable performance across a range of multi-modal tasks, they face the persistent challenge of hallucination, which introduces practical weaknesses and raises concerns about their reliable deployment in real-world applications. Existing work has explored contrastive decoding approaches to mitigate this issue, where the output of the original LVLM is compared and contrasted with that of a perturbed version. However, these methods require two or more queries that slow down LVLM response generation, making them less suitable for real-time applications. To overcome this limitation, we propose ONLY, a training-free decoding approach that requires only a single query and a one-layer intervention during decoding, enabling efficient real-time deployment. Specifically, we enhance textual outputs by selectively amplifying crucial textual information using a text-to-visual entropy ratio for each token. Extensive experimental results demonstrate that our proposed ONLY consistently outperforms state-of-the-art methods across various benchmarks while requiring minimal implementation effort and computational cost. Code is available at <https://github.com/zifuwang/ONLY>.

## 1. Introduction

Recent advances in large vision-language models (LVLMs), which expand the capabilities of large language models (LLMs) to visual understanding and reasoning [1, 17, 37], have demonstrated exceptional performance across various vision-language tasks, such as object detection [36, 40], segmentation [15, 34], and image captioning [20, 28]. However, a persistent challenge with current LVLMs is their tendency to generate hallucinated content, where the generated responses do not align accurately with the actual image input [24]. This can significantly impact the reliability

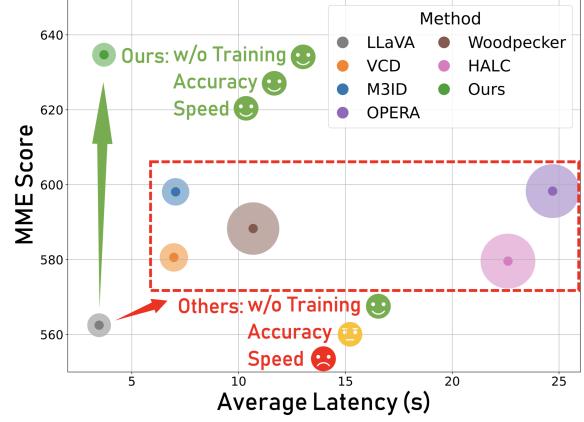


Figure 1. Comparisons of accuracy and inference speed of multiple hallucination mitigation approaches. The size of bubbles stands for the GPU memory consumption. Our method effectively mitigates hallucination with only 0.07× extra time.

of LVLMs in real-world applications where precise visual interpretation is essential [2, 3, 24]. Therefore, addressing hallucinations in LVLMs is crucial to ensuring their safe and effective deployment in critical domains.

To alleviate the hallucination problem, early work identified biased training sets as a critical cause and, as a result, attempted to establish curated training datasets and adopt robust fine-tuning techniques [5, 31, 44]. However, their reliance on additional data and the need for fine-tuning large-scale models make these approaches time-consuming and impractical for individual users. Another common approach is contrastive decoding [18], which eliminates the need for costly training by directly intervening in the decoding process during inference. Specifically, these methods typically introduce a distorted set of inputs, and contrast their respective token predictions with the predictions from original data to mitigate undesired hallucinations [8, 10, 16, 35]. Although existing contrastive decoding-based approaches achieve notable performance improvements, they require multiple LVLM queries to process both the original and distorted inputs, resulting in response times that are twice as long, or more, making

\*Equal contribution. Contact: zifuw@cs.cmu.edu.

them less suitable for real-time applications [4, 10, 16].

To illustrate this, we analyze the performance-efficiency trade-off of existing approaches for mitigating hallucinations in LVLMs and present the results in Figure 1. As we can see, while other hallucination mitigation methods achieve higher accuracy on the hallucination evaluation benchmark, they come at a significant cost, requiring 2× or more inference time and higher GPU memory consumption. We recognize this overhead is impractical given the limited performance improvements, highlighting the urgent need for more efficient approaches that can effectively mitigate hallucinations in LVLMs.

In this work, we introduce ONLY, a training-free approach that requires only a single query and a one-layer intervention during decoding, offering an efficient solution for mitigating hallucinations in LVLMs. Our ONLY approach selects attention heads that prioritize textual information over visual information—specifically, those with a high text-to-visual entropy ratio—to stimulate textually enhanced next-token predictions. The enhanced output is then adaptively contrasted/collaborated with the original output logits using a single-layer intervention, aiming to reduce predominant and irrelevant language bias. Our ONLY approach is both simple and effective, requiring just one additional attention layer computation. It incurs a modest 1.07× increase in inference time with negligible GPU memory overhead, significantly lower than the 2× or more increase seen in previous contrastive decoding methods. Moreover, ONLY achieves superior performance across multiple benchmarks, outperforming the current state-of-the-art by 3.14% on POPE and 1.6% on CHAIR.

To validate the effectiveness of our proposed ONLY approach, we evaluate it on three LVLMs (*i.e.*, LLaVA-1.5 [22], InstructBLIP [9], and Qwen-VL [1]) across various benchmarks, including POPE [19], CHAIR [28], MME-Hallucination [11], MMBench [25], MM-Vet [39], and MMVP [32]. Extensive experimental results demonstrate that our ONLY approach consistently outperforms state-of-the-art methods across these benchmarks while requiring minimal implementation effort and computational cost. Additionally, qualitative case studies and GPT-4V-aided evaluations on LLaVA-Bench further validate the effectiveness of our ONLY approach in enhancing the coherence and accuracy of LVLM responses.

Our contributions are summarized as follows:

- We investigate and challenge the performance-efficiency trade-off of existing contrastive decoding approaches for mitigating hallucinations in LVLMs, highlighting the efficiency issues.
- We present ONLY, a novel training-free decoding algorithm that leverages a single additional Transformer layer to improve the accuracy of LVLM responses.
- We conduct comprehensive experiments across various

benchmarks and demonstrate that our proposed ONLY consistently outperforms existing approaches with minimal implementation effort and computational cost.

## 2. Related Work

**Large Vision-Language Models (LVLMs).** Recently, large-scale LLMs have demonstrated remarkable proficiency in handling human queries and exhibit robust linguistic capabilities [7, 33]. Leveraging these powerful models, researchers are exploring ways to align the visual modality with language, unlocking advanced visual recognition and reasoning capabilities across various multi-modal tasks [2, 24]. For example, LLaVA-1.5 [21] employs a pre-trained CLIP ViT-L/14 [27] as the vision encoder, and trains a linear mapping layer to connect the vision and language modalities. In contrast, InstructBLIP [9] builds on a pre-trained BLIP-2 [17] and incorporates an instruction-aware Q-Former module to bridge the modalities. Despite their exceptional multi-modal performance, these LVLMs still suffer from hallucinations, often generating text responses that do not accurately reflect the given image input [3, 26, 31, 41, 43]. Such hallucinations pose significant challenges for deploying these models in real-world applications. In this work, we propose a novel training-free algorithm that mitigates hallucinations while improving the efficiency of LVLMs for real-world deployment.

**Hallucination in LVLMs.** Recent studies have revealed that LVLMs may generate cross-modal inconsistencies between visual inputs and their corresponding responses, *i.e.*, hallucinations, which can lead to misinformation and performance degradation [16, 24]. To mitigate these hallucinations, early works have explored the use of additional robust instruction tuning on curated datasets [14, 26, 31]. While effective, these methods require extensive and costly training, making them impractical for individual users. More recently, researchers have explored an alternative approach by developing variant methods based on contrastive decoding strategies, which mitigate hallucinations and enhance coherence by contrasting logits from counterpart outputs [6, 10, 16]. However, we recognize that these methods require two or even multiple queries, which slows down LVLM response generation, making them less suitable for real-time applications. In response, we propose ONLY, a contrastive decoding-based approach that requires only a one-time query and a one-layer intervention during decoding, achieving competitive performance while effectively minimizing implementation efforts and computational costs.

## 3. Method

In this work, we present ONLY, a training-free algorithm that uses only one Transformer layer to improve the accuracy of LVLM responses, as illustrated in Figure 2.

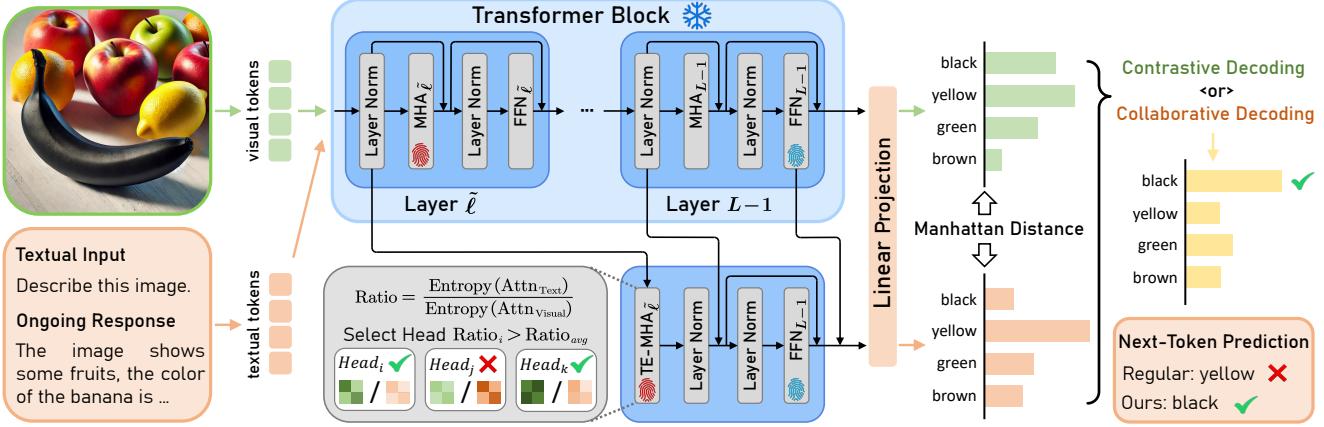


Figure 2. **Overview of our proposed ONLY.** Our method retains the core decoding process of LVLMs but incorporates a textual-enhanced multi-head attention layer with a residual connection to the last layer’s output. This adjustment aims to produce an output with a greater focus on textual information. The resulting textual-enhanced logits are then adaptively decoded alongside the original output, employing either contrastive or collaborative decoding strategies to optimize performance.

### 3.1. Preliminaries

**LVLM Decoding.** Recent LVLMs effectively process both visual and linguistic data using three key components: vision encoders, connectors, and a Large Language Model (LLM). An LVLM, parameterized by  $\theta$ , autoregressively generates a fluent textual response sequence  $y$  from an input image  $v$  and a textual query  $x$ . Initially,  $v$  is processed by a vision encoder and transformed into visual tokens via a vision-language alignment module (*e.g.*, Q-Former [17] or linear projection [21]). These tokens, combined with the query tokens, are input to the LLM. The generation of each token  $y_t$  in the sequence  $y$  is modeled as:

$$y_t \sim p_\theta(y_t|v, x, y_{<t}) = \text{softmax}(f_\theta(y_t|v, x, y_{<t}))_{y_t}, \quad (1)$$

where  $y_t \in \mathcal{S}$  is the current token,  $y_{<t} = [y_0, \dots, y_{t-1}]$  are the previously generated tokens, and  $f_\theta$  represents the logits over a vocabulary set  $\mathcal{S}$ .

**Transformer Decoder.** The language model is structured as a Transformer, where a sequence of tokens  $\{x_1, x_2, \dots, x_{t-1}\}$  are initially embedded into a sequence of hidden states  $\mathcal{H}_{t-1}^0 = \{h_1^0, \dots, h_{t-1}^0\}$ . The Transformer comprises  $L$  layers, each layer incorporates a Multi-Head Attention (MHA) module and a Multi-Layer Perceptron (MLP). At time step  $t$ , the output of each layer  $\mathcal{H}_t^{\ell+1}$  is derived from the hidden states input  $\mathcal{H}_t^\ell$ , employing two primary residual connections:

$$\bar{\mathcal{H}}_t^\ell = \text{MHA}_\ell(\mathcal{H}_t^\ell) + \mathcal{H}_t^\ell, \quad \mathcal{H}_t^{\ell+1} = \text{MLP}_\ell(\bar{\mathcal{H}}_t^\ell) + \bar{\mathcal{H}}_t^\ell. \quad (2)$$

Each MHA module consists of  $H$  attention heads that compute self-attention, where the attention score is derived from query, key, and value matrices. Specifically, for the  $i$ -th

head in layer  $\ell$ , the operation is given by:

$$\begin{aligned} \text{Head}_{\ell,i}(\mathcal{H}_t^\ell) &= \text{Attention}(Q_{\ell,i}, K_{\ell,i}, V_{\ell,i}) \\ &= \text{softmax}\left(\frac{Q_{\ell,i} \cdot K_{\ell,i}^\top}{\sqrt{d_k}}\right) V_{\ell,i}, \end{aligned} \quad (3)$$

where  $Q_{\ell,i}/K_{\ell,i}/V_{\ell,i} = \mathcal{H}_t^\ell W_{\ell,i}^{Q/K/V}$ , are query/key/value matrices obtained from learned weights. The outputs from all  $H$  heads are then concatenated and projected using an projection matrix  $W_\ell^O$ :

$$\begin{aligned} \text{MHA}_\ell(\mathcal{H}_t^\ell) &= \text{Concat}(\text{Head}_{\ell,1}(\mathcal{H}_t^\ell), \\ &\quad \text{Head}_{\ell,2}(\mathcal{H}_t^\ell), \dots, \text{Head}_{\ell,H}(\mathcal{H}_t^\ell)) W_\ell^O. \end{aligned} \quad (4)$$

At last, a projection head  $\phi(\cdot)$  predicts the logits of the next token  $x_t$  over the vocabulary set  $\mathcal{S}$ :

$$f_\theta(y_t|y_{<t}) = \phi(\mathcal{H}_t^L), y_t \in \mathcal{S}. \quad (5)$$

Combining Eq. 5 with Eq. 1, we finally obtain:

$$p_\theta(y_t|v, x, y_{<t}) = \text{softmax}(\phi(\mathcal{H}_t^L))_{y_t}. \quad (6)$$

### 3.2. One-Layer Intervention for Textual Enhancement

Previous contrastive decoding methods focus primarily on the visual modality or the effect of visual input on the textual modality: *e.g.*, VCD [16] contrasts the outputs obtained with original vs. visual distorted input, and M3ID [10] amplifies the influence of the reference image over the language prior. However, the effect of textual modality has been less studied. In this work, we propose to address hallucination by directly producing textually-enhanced outputs with minimal additional computational

---

**Algorithm 1** Predict Textual-Enhanced (TE) Logits

---

**Require:** Initial hidden states  $\mathcal{H}_t^0$ , total transformer layers  $L$ , total attention heads  $H$ , layer index for textual enhancement  $\tilde{\ell}$ .

- 1: **procedure** PREDICT\_TE\_LOGITS( $A$ )
- 2:   **for**  $\ell \in \{0, 1, 2, \dots, L-1\}$  **do**
- 3:     **for**  $i \in \{0, 1, \dots, H-1\}$  **do**
- 4:       **Step 1: Calculate TE attention output**
- 5:        **if**  $\ell = \tilde{\ell}$  **then**
- 6:           $\tilde{\mathcal{H}}_t^{\tilde{\ell}} \leftarrow \text{TE-MHA}_{\tilde{\ell}}(\mathcal{H}_t^{\tilde{\ell}})$        $\triangleright$  Equation 13
- 7:        **end if**
- 8:     **end for**
- 9:     **Step 2: Calculate Transformer output for each layer**
- 10:     $\mathcal{H}_t^{\ell} \leftarrow \text{MHA}_{\ell}(\mathcal{H}_t^{\ell}) + \tilde{\mathcal{H}}_t^{\tilde{\ell}}$        $\triangleright$  Equation 2
- 11:     $\mathcal{H}_t^{\ell+1} \leftarrow \text{MLP}_{\ell}(\mathcal{H}_t^{\ell}) + \mathcal{H}_t^{\ell}$        $\triangleright$  Equation 2
- 12:   **if**  $\ell = L-1$  **then**
- 13:     **Step 3: Calculate TE Transformer output**
- 14:      $\tilde{\mathcal{H}}_t^{L-1} \leftarrow \tilde{\mathcal{H}}_t^{\tilde{\ell}} + \mathcal{H}_t^{L-1}$        $\triangleright$  Equation 15
- 15:      $\hat{\mathcal{H}}_t^L \leftarrow \text{MLP}_{L-1}(\tilde{\mathcal{H}}_t^{L-1}) + \tilde{\mathcal{H}}_t^{L-1}$        $\triangleright$  Equation 16
- 16:   **end if**
- 17:   **end for**
- 18:   **Step 4: Calculate original logits and TE logits**
- 19:    Logits =  $f_{\theta}(y_t | v, \mathbf{x}, \mathbf{y}_{<t}) \leftarrow \text{Linear}(\mathcal{H}_t^L)$
- 20:    Logits\_TE =  $f_{\theta}(y_t | v, \mathbf{x}, \mathbf{y}_{<t}) \leftarrow \text{Linear}(\hat{\mathcal{H}}_t^L + \mathcal{H}_t^L)$
- 21:   **return** Logits, Logits\_TE
- 22: **end procedure**

---

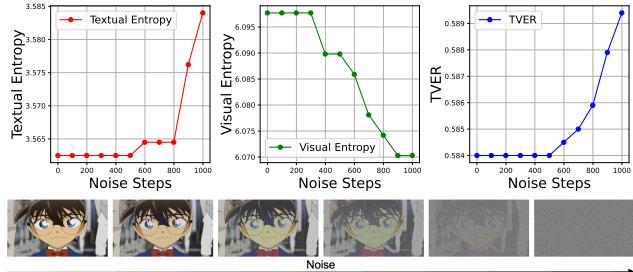


Figure 3. **Impact of applying diffusion noise on textual and visual attention entropy.** We perform an analysis on all COCO samples from the POPE benchmark and observe that as distortion increases, textual entropy rises whereas visual entropy decreases.

overhead. Specifically, inspired by information theory [30], we introduce an attention-head selection strategy guided by the text-to-visual entropy ratio. As illustrated in Figure 3, we observe that when the distortion escalates (similar to the diffusion steps in VCD), textual entropy increases while visual entropy decreases. Guided by this observation, we propose to directly select attention heads with a higher text-to-visual entropy ratio to stimulate textually-enhanced outputs to avoid double queries while extracting language bias.

**Attention Head Selection Using Text-to-Visual Entropy Ratio.** Suppose a token is generated at time step  $t$ , and the initial hidden states input to the Transformer decoder for this token is  $\mathcal{H}_t^0$ . For layer  $\ell$ , the input hidden states can be denoted as  $\mathcal{H}_t^{\ell}$ . We distinguish between text

and visual attention within the attention matrix by computing the raw attention scores  $a_{\ell,i}$  for each head:

$$a_{\ell,i} = \text{softmax}(Q_{\ell,i} \cdot K_{\ell,i}^{\top} / \sqrt{d_k}), \quad (7)$$

where  $Q_{\ell,i}$  and  $K_{\ell,i}$  are the query and key matrices for head  $i$  in layer  $\ell$ . To isolate text and visual attentions, we utilize indices corresponding to textual or visual tokens:

$$a_{\ell,i}^T = \{a_{\ell,i,j} \mid j \in \text{indices}_T\}, a_{\ell,i}^V = \{a_{\ell,i,j} \mid j \in \text{indices}_V\}, \quad (8)$$

where  $\text{indices}_T$  and  $\text{indices}_V$  specify positions of textual and visual tokens, respectively. The entropy for these attention sets is computed as follows:

$$\begin{aligned} \text{Entropy}(a_{\ell,i}^T) &= - \sum_k p_{\ell,i,k}^T \log p_{\ell,i,k}^T, \\ \text{Entropy}(a_{\ell,i}^V) &= - \sum_k p_{\ell,i,k}^V \log p_{\ell,i,k}^V, \end{aligned} \quad (9)$$

where  $p_{\ell,i,k}^T$  and  $p_{\ell,i,k}^V$  represent the normalized attention probabilities, computed from the softmax of each subset:

$$p_{\ell,i,k}^T = \text{softmax}(a_{\ell,i,k}^T), p_{\ell,i,k}^V = \text{softmax}(a_{\ell,i,k}^V). \quad (10)$$

The Text-to-Visual Entropy Ratio (TVER) for each attention head is calculated as:

$$\text{TVER}_{\ell,i} = \frac{\text{Entropy}(a_{\ell,i}^T)}{\text{Entropy}(a_{\ell,i}^V)}. \quad (11)$$

To optimize the attention output for enhanced textual relevance while reducing visual information, we selectively deactivate heads with a TVER below the average for that layer, setting their attention weights to zero. This approach prioritizes heads with relatively higher text-to-visual entropy ratios, providing a clue where uncertainty in the textual modality is higher:

$$\tilde{a}_{\ell,i} = \begin{cases} a_{\ell,i}, & \text{if } \text{TVER}_{\ell,i} \geq \text{average}(\text{TVER}_{\ell}), \\ 0, & \text{otherwise.} \end{cases} \quad (12)$$

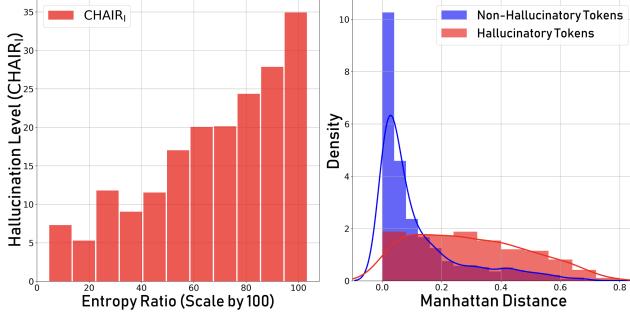
With this, we obtain the output of the Textual-Enhanced Multi-Head Attention (TE-MHA) module:

$$\begin{aligned} \text{TE-MHA}_{\ell}(\mathcal{H}_t^{\ell}) &= \\ \text{Concat}(\tilde{a}_{\ell,1}V_{l,1}, \tilde{a}_{\ell,2}V_{l,2}, \dots, \tilde{a}_{\ell,H}V_{l,H})W_{\ell}^O. \end{aligned} \quad (13)$$

### 3.3. Adaptive Decoding

In this section, we utilize the logits obtained from textual-enhanced attention outputs for adaptive decoding.

Suppose layer  $\tilde{\ell} \in \{0, 1, \dots, L-1\}$  is the selected layer for textual enhancement, where we calculate a textual-enhanced attention output as discussed in Eq. 13. To ensure



**Figure 4. Text-to-visual entropy ratio is correlated with hallucinations.** (Left) Density plot of token-wise average textual-to-visual entropy ratio and bar plot of average  $\text{CHAIR}_I$  in each bin on the CHAIR benchmark; (Right) Density plots of token-level Manhattan distance between original and textual-enhanced logits for both hallucinatory and non-hallucinatory tokens on POPE.

that the output logits do not deviate excessively from the original LVLM outputs, we implement two residual connections. These connections are defined as follows:

$$\tilde{\mathcal{H}}_t^{\ell} = \text{TE-MHA}_{\ell}(\mathcal{H}_t^{\ell}), \quad (14)$$

$$\tilde{\mathcal{H}}_t^{L-1} = \tilde{\mathcal{H}}_t^{\ell} + \mathcal{H}_t^{L-1}, \quad (15)$$

$$\hat{\mathcal{H}}_t^L = \text{MLP}_{L-1}(\tilde{\mathcal{H}}_t^{L-1}) + \tilde{\mathcal{H}}_t^{L-1}. \quad (16)$$

Finally, the textual-enhanced predicted probability can be obtained by:

$$\begin{aligned} \tilde{p}_{\theta}(y_t|v, \mathbf{x}, \mathbf{y}_{<t}) &= \text{softmax}(\tilde{f}_{\theta}(y_t|v, \mathbf{x}, \mathbf{y}_{<t}))_{y_t} \\ &= \text{softmax}(\phi(\hat{\mathcal{H}}_t^L))_{y_t}. \end{aligned} \quad (17)$$

To adaptively contrast the original and textual-enhanced logits, we measure the Manhattan distance between the two probability distributions at each timestep  $t$ :

$$d_t = \sum_{y_t \in \mathcal{S}} |p_{\theta}(y_t|v, \mathbf{x}, \mathbf{y}_{<t}) - \tilde{p}_{\theta}(y_t|v, \mathbf{x}, \mathbf{y}_{<t})|, \quad (18)$$

where  $d_t$  provides a measure of the difference between the distributions. Based on this distance, we adjust the original logits either collaboratively or contrastively:

$$y_t \sim p_{\theta}(y_t) = \text{softmax}(f_{\theta}^{\text{final}}) \quad (19)$$

$$f_{\theta}^{\text{final}} = \begin{cases} f_{\theta}(y_t|v, \mathbf{x}, \mathbf{y}_{<t}) + \alpha_1 \tilde{f}_{\theta}(y_t|v, \mathbf{x}, \mathbf{y}_{<t}), & \text{if } d_t < \gamma \text{ (collaborative);} \\ (1 + \alpha_2) f_{\theta}(y_t|v, \mathbf{x}, \mathbf{y}_{<t}) - \alpha_2 \tilde{f}_{\theta}(y_t|v, \mathbf{x}, \mathbf{y}_{<t}), & \text{if } d_t \geq \gamma \text{ (contrastive),} \end{cases} \quad (20)$$

where  $\gamma$  is a predefined threshold that determines the decoding strategy based on the measured distance.

**Effectiveness of text-to-visual entropy ratio for textual information enhancement.** We further conduct an empirical study to validate the effectiveness of applying text-to-visual entropy ratio for language bias reflection, as shown in Figure 4. The experimental results demonstrate that *the entropy ratio is strongly correlated to the hallucination level at both the response and token levels*.

## 4. Experiments

In this section, we evaluate the effectiveness of our method in mitigating hallucinations in LVLMs across a range of benchmarking scenarios, comparing it with existing state-of-the-art approaches.

### 4.1. Experimental Settings

**Evaluated LVLMs.** We evaluate the effectiveness of our method on three state-of-the-art open-source LVLMs: LLaVA-1.5 [22], InstructBLIP [9] and Qwen-VL [1].

**Benchmarks.** We conduct extensive experiments on six benchmarks: (1) **POPE** [19] is a benchmark commonly used to assess object hallucinations in LVLMs, which evaluates model accuracy through yes-or-no questions about the presence of specific objects in images; (2) **CHAIR** [28] evaluates object hallucinations through image captioning, where the LVLMs are prompted to describe 500 randomly selected images from the MSCOCO validation set; (3) **MME-Hallucination** [11] is a comprehensive benchmark for LVLMs consisting of four subsets: *existence* and *count* for object-level hallucinations, and *position* and *color* for attribute-level hallucinations; (4) **MMBench** [25] is a benchmark for evaluating LVLMs' multi-modal understanding ability across 20 dimensions; (5) **MMVP** [32] comprises 150 CLIP-blind image pairs, each paired with a binary-option question to evaluate the fine-grained visual recognition capabilities of LVLMs; (6) **MM-Vet** [39] utilizes LLM-based evaluator to evaluate LVLMs on 6 capabilities, including recognition, OCR, knowledge, language generation, spatial awareness, and math; (7) **LLaVA-Bench** provides 24 images in complex scenes, memes, and sketches, along with 60 challenging questions.

**Baselines.** We compare the performance of our ONLY approach with the following state-of-the-art approaches: VCD [16], M3ID [10], Woodpecker [38], HALC [6], DoLa [8] and OPERA [12]. We apply sampling-based decoding in default, where the next token is sampled directly from the post-softmax probability distribution.

**Implementation Details.** We follow the default query format for all LVLMs. Besides, we set  $\alpha_1 = 3$ ,  $\alpha_2 = 1$ , and  $\gamma = 0.2$  for LLaVA-1.5 [22], and  $\gamma = 0.4$  for InstructBLIP [9] / Qwen-VL [1]. Following VCD [16], we implement adaptive plausibility constraints [18] with  $\beta = 0.1$  across all tasks. All experiments are performed on a single 48GB NVIDIA RTX 6000 Ada GPU.

	Setup	Method	LLaVA-1.5				InstructBLIP				Qwen-VL			
			Acc. $\uparrow$	Prec. $\uparrow$	Rec. $\uparrow$	F1 $\uparrow$	Acc. $\uparrow$	Prec. $\uparrow$	Rec. $\uparrow$	F1 $\uparrow$	Acc. $\uparrow$	Prec. $\uparrow$	Rec. $\uparrow$	F1 $\uparrow$
MS-COCO	Random	Regular	83.13	81.94	85.00	83.44	83.07	83.02	83.13	83.08	85.23	97.23	72.53	83.09
		VCD	87.00	86.13	<u>88.20</u>	87.15	86.23	<u>88.14</u>	83.73	85.88	<u>87.03</u>	97.36	<u>76.13</u>	<u>85.45</u>
		M3ID	<u>87.50</u>	<u>87.38</u>	87.67	<u>87.52</u>	<u>86.67</u>	88.09	<u>84.80</u>	<u>86.41</u>	86.40	<u>98.23</u>	74.13	84.50
		<b>Ours</b>	<b>89.70</b>	<b>89.95</b>	<b>88.27</b>	<b>89.10</b>	<b>89.23</b>	<b>91.83</b>	<b>86.13</b>	<b>88.89</b>	<b>88.90</b>	<b>98.52</b>	<b>79.27</b>	<b>87.85</b>
A-OKVQA	Popular	Regular	81.17	78.28	86.27	82.08	77.00	73.82	83.67	78.44	84.53	94.50	73.33	82.58
		VCD	83.10	79.96	<u>88.33</u>	83.94	80.07	77.67	84.40	80.89	85.87	94.98	<u>75.73</u>	84.27
		M3ID	84.30	<u>81.58</u>	<b>88.60</b>	<u>84.95</u>	<u>80.97</u>	<u>77.93</u>	<b>86.40</b>	<u>81.85</u>	<u>86.07</u>	<b>96.56</b>	74.80	<u>84.30</u>
		<b>Ours</b>	<b>86.00</b>	<b>84.44</b>	88.27	<b>86.31</b>	<b>83.27</b>	<b>81.46</b>	<u>86.13</u>	<b>83.73</b>	<b>87.47</b>	<u>95.63</u>	<b>79.48</b>	<b>86.81</b>
GQA	Adversarial	Regular	77.43	73.31	86.27	79.26	74.60	71.26	82.47	76.45	83.37	91.47	73.60	81.57
		VCD	77.17	72.18	<b>88.40</b>	79.47	77.20	<u>74.29</u>	83.20	78.49	<u>83.73</u>	89.84	<u>76.07</u>	<u>82.38</u>
		M3ID	<u>78.23</u>	<u>73.51</u>	88.27	<u>80.22</u>	<u>77.47</u>	73.68	<u>85.47</u>	79.14	83.37	91.19	73.87	81.62
		<b>Ours</b>	<b>79.40</b>	<b>75.00</b>	88.20	<b>81.07</b>	<b>80.10</b>	<b>76.89</b>	<b>86.07</b>	<b>81.22</b>	<b>83.80</b>	<b>92.33</b>	<b>76.14</b>	<b>83.46</b>
MS-COCO	Random	Regular	81.90	76.63	91.80	83.53	80.63	76.82	87.73	81.92	86.40	94.32	77.47	85.07
		VCD	83.83	78.05	<u>94.13</u>	85.34	84.20	80.90	89.53	85.00	<u>87.93</u>	94.59	<u>80.47</u>	<u>86.96</u>
		M3ID	<u>84.67</u>	<u>79.25</u>	93.93	<u>85.97</u>	<u>85.43</u>	<u>81.77</u>	<u>91.20</u>	<u>86.23</u>	87.50	<u>95.33</u>	78.87	86.32
		<b>Ours</b>	<b>86.07</b>	<b>80.91</b>	<b>94.40</b>	<b>87.14</b>	<b>88.57</b>	<b>86.13</b>	<b>91.93</b>	<b>88.94</b>	<b>89.47</b>	<b>95.34</b>	<b>83.84</b>	<b>89.22</b>
A-OKVQA	Popular	Regular	75.07	68.58	92.53	78.77	75.17	70.15	87.60	77.91	85.77	92.82	77.53	84.49
		VCD	76.63	69.59	<b>94.60</b>	80.19	78.63	<u>73.53</u>	89.47	80.72	87.33	93.68	<u>80.07</u>	<u>86.34</u>
		M3ID	<u>77.80</u>	<u>70.98</u>	94.07	<u>80.91</u>	<u>78.80</u>	73.38	<u>90.40</u>	81.00	87.37	<b>95.31</b>	78.60	86.15
		<b>Ours</b>	<b>79.00</b>	<b>72.17</b>	<u>94.40</u>	<b>81.80</b>	<b>80.83</b>	<b>75.23</b>	<b>91.93</b>	<b>82.75</b>	<b>89.47</b>	<u>94.77</u>	<b>84.43</b>	<b>89.30</b>
GQA	Adversarial	Regular	67.23	61.56	91.80	73.70	69.87	64.54	88.20	74.54	80.37	82.56	77.00	79.68
		VCD	67.40	61.39	93.80	74.21	<u>71.00</u>	<u>65.41</u>	89.13	<u>75.45</u>	<u>81.90</u>	83.07	<u>80.13</u>	<u>81.57</u>
		M3ID	<u>68.60</u>	<u>62.22</u>	<b>94.73</b>	<u>75.11</u>	70.10	64.28	<u>90.47</u>	75.16	<u>81.90</u>	84.25	78.47	81.26
		<b>Ours</b>	<b>68.70</b>	<b>62.35</b>	<u>94.40</u>	<b>75.70</b>	<b>72.47</b>	<b>66.19</b>	<b>91.87</b>	<b>76.94</b>	<b>82.07</b>	<b>85.02</b>	<b>81.09</b>	<b>83.01</b>
MS-COCO	Random	Regular	82.23	76.32	93.47	84.03	79.67	76.05	86.60	80.99	85.10	91.42	77.47	83.87
		VCD	83.23	76.73	<u>95.40</u>	85.05	82.83	<u>80.16</u>	87.27	83.56	87.00	92.11	<u>80.93</u>	<u>86.16</u>
		M3ID	84.20	78.00	95.27	<u>85.77</u>	83.07	80.06	88.07	83.87	87.07	92.64	80.53	<u>86.16</u>
		<b>Ours</b>	<b>86.70</b>	<b>80.94</b>	<b>96.00</b>	<b>87.83</b>	<b>86.17</b>	<b>83.84</b>	<b>89.60</b>	<b>86.63</b>	<b>88.03</b>	<b>93.59</b>	<b>82.68</b>	<b>87.80</b>
A-OKVQA	Popular	Regular	73.47	<b>66.83</b>	93.20	77.84	73.33	68.72	85.67	76.26	80.87	82.65	78.13	80.33
		VCD	72.37	65.27	<u>95.60</u>	77.58	<u>76.13</u>	<u>71.10</u>	88.07	<u>78.68</u>	82.53	83.52	<u>81.07</u>	<u>82.27</u>
		M3ID	73.87	<u>66.70</u>	95.33	<u>78.49</u>	75.17	69.94	88.27	78.04	82.68	83.74	80.85	<u>82.27</u>
		<b>Ours</b>	<b>74.03</b>	<u>66.70</u>	<b>96.00</b>	<b>78.71</b>	<b>77.20</b>	<b>71.79</b>	<b>89.60</b>	<b>79.72</b>	<b>82.87</b>	<b>83.88</b>	<b>82.55</b>	<b>83.21</b>
GQA	Adversarial	Regular	68.60	<u>62.43</u>	93.40	74.84	68.60	63.94	85.33	73.10	78.77	79.33	77.80	78.56
		VCD	<u>68.83</u>	62.26	<u>95.67</u>	<u>75.43</u>	71.00	65.75	87.67	75.14	81.17	81.48	<u>80.67</u>	81.07
		M3ID	68.67	62.16	95.40	75.28	71.17	<u>65.79</u>	<b>88.20</b>	<u>75.36</u>	<b>81.90</b>	<b>83.07</b>	80.13	<u>81.57</u>
		<b>Ours</b>	<b>69.23</b>	<b>62.55</b>	<b>95.87</b>	<b>75.70</b>	<b>71.93</b>	<b>65.98</b>	<u>87.93</u>	<b>75.84</b>	<u>81.33</u>	<u>82.38</u>	<b>81.50</b>	<b>81.94</b>

Table 1. **Results on POPE [19] benchmark.** Higher ( $\uparrow$ ) accuracy, precision, recall, and F1 indicate better performance. The best results are **bolded**, and the second-best are underlined.

Method	LLaVA-1.5				InstructBLIP				Qwen-VL			
	Max Token 64		Max Token 128		Max Token 64		Max Token 128		Max Token 64		Max Token 128	
	CHAIR <sub>S</sub> $\downarrow$	CHAIR <sub>I</sub> $\downarrow$										
Regular	26.2	9.4	55.0	16.3	31.2	11.1	57.0	17.6	33.6	12.9	52.0	16.5
VCD	24.4	7.9	54.4	16.6	30.0	10.1	60.4	17.8	33.0	12.8	50.2	16.8
M3ID	<u>21.4</u>	<u>6.3</u>	56.6	15.7	30.8	10.4	62.2	18.1	32.2	11.5	<u>49.5</u>	17.2
Woodpecker	24.9	7.5	57.6	16.7	31.2	10.8	60.8	17.6	31.1	12.3	51.8	16.3
HALC	21.7	7.1	<u>51.0</u>	<u>14.8</u>	<u>24.5</u>	<b>8.0</b>	<u>53.8</u>	<u>15.7</u>	<u>28.2</u>	<u>9.1</u>	49.6	<u>15.4</u>
<b>Ours</b>	<b>20.0</b>	<b>6.2</b>	<b>49.8</b>	<b>14.3</b>	<b>23.5</b>	<u>8.2</u>	<b>52.2</b>	<b>15.5</b>	<b>27.3</b>	<b>8.4</b>	<b>48.0</b>	<b>14.3</b>

Table 2. **Results on CHAIR [28] benchmark.** We limit the maximum number of new tokens to 64 or 128. Lower ( $\downarrow$ ) CHAIR<sub>S</sub>, CHAIR<sub>I</sub> indicate better performance. The best results in each setting are **bolded**, and the second-best are underlined.

Method	Object-level		Attribute-level		MME Score $\uparrow$
	Existence $\uparrow$	Count $\uparrow$	Position $\uparrow$	Color $\uparrow$	
Regular	173.75	121.67	117.92	149.17	562.50
DoLa	176.67	113.33	90.55	141.67	522.22
OPERA	183.33	137.22	122.78	155.00	598.33
VCD	186.67	125.56	128.89	139.45	580.56
M3ID	186.67	128.33	131.67	151.67	598.11
Woodpecker	<u>187.50</u>	125.00	126.66	149.17	588.33
HALC	183.33	133.33	107.92	<u>155.00</u>	579.58
<b>Ours</b>	<b>191.67</b>	<b>145.55</b>	<b>136.66</b>	<b>161.66</b>	<b>635.55</b>

Table 3. **Results on MME-Hallucination [11] with LLaVA-1.5 [22]**. We report the average MME scores for each subset. Higher scores ( $\uparrow$ ) indicate better performance. The best results are **bolded**, and the second-best are underlined.

## 4.2. Results and Discussions

**Results on POPE.** In Table 1, we compare our method’s performance against various baselines on the POPE benchmark. As shown in the table, our approach consistently outperforms previous state-of-the-art methods across various LVLM models and settings, demonstrating its robustness across different evaluation scenarios. Specifically, in the MS-COCO (Random) setting with the LLaVA-1.5 backbone, our method surpasses VCD by 2.20% and M3ID by 1.70% in accuracy. Even in the more challenging adversarial setting, our approach maintains its superior performance, outperforming VCD by 2.23% and M3ID by 1.17%. Overall, these consistent gains across different datasets and LVLM models highlight the effectiveness of our method as a strong and generalizable solution for mitigating hallucinations in LVLMs.

**Results on CHAIR.** On the open-ended CHAIR benchmark, our ONLY method achieves superior performance with lower hallucination rates. Table 2 presents a comparison against four state-of-the-art approaches, evaluating hallucination rates with  $\text{CHAIR}_S$  and  $\text{CHAIR}_I$  under maximum token generation limits of 64 and 128 across three LVLM backbones. Notably, in the LLaVA-1.5 (Max Token = 128) setting, our approach reduces  $\text{CHAIR}_S$  by 5.2 points and  $\text{CHAIR}_I$  by 2.0 points compared to regular decoding.

**Results on MME.** In Table 3, we compare our approach against other methods on the MME benchmark. The results show that our method consistently outperforms all baselines, achieving the highest scores across both object-level (Existence, Count) and attribute-level (Position, Color) evaluations. Notably, our method attains an MME score of 634.67, outperforming the second-best method, M3ID, by 36.34 points, demonstrating its superior capability in mitigating various types of hallucinations.

**Results on MMVP.** To evaluate the effectiveness of our approach on fine-grained visual recognition tasks, we conduct experiments on the MMVP benchmark and present the results in Figure 5. With our ONLY approach, the LVLM is able to handle more nuanced visual recognition tasks, improving the performance from 22.67% to 28.00%.

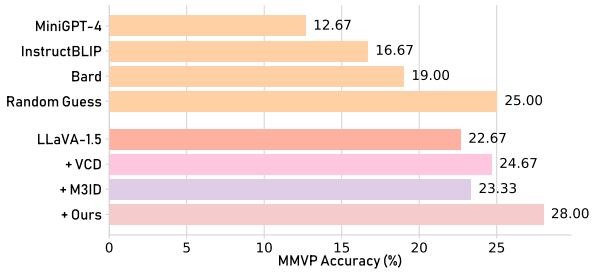


Figure 5. **Results on MMVP [32]**. We apply our approach to LLaVA-1.5 [22] and compare its performance against other hallucination mitigation methods.

Method	LLaVA-1.5		InstructBLIP	
	Acc. $\uparrow$	Det. $\uparrow$	Acc. $\uparrow$	Det. $\uparrow$
Regular	6.07	6.20	5.26	5.53
<b>Ours</b>	<b>7.00</b>	<b>7.13</b>	<b>6.60</b>	<b>6.73</b>
VCD	4.60	5.13	4.87	5.33
<b>Ours</b>	<b>6.27</b>	<b>6.60</b>	<b>6.80</b>	<b>6.93</b>
M3ID	6.13	6.27	5.93	6.20
<b>Ours</b>	<b>6.73</b>	<b>6.80</b>	<b>6.67</b>	<b>6.87</b>

Table 4. **GPT-4V-aided evaluation on LLaVA-Bench**. Higher accuracy and detailedness ( $\uparrow$ ) indicate better performance. The evaluation is performed on LLaVA-1.5 [22].

**Results on MMBench and MMVet.** We also report the performance of all compared methods on the MMBench and MMVet benchmarks in Table 5. Our approach continues to outperform existing state-of-the-art methods, demonstrating that it also enhances the general multi-modal understanding capabilities of LVLMs.

**Results on LLaVA-Bench.** In Figure 6, we present a case study on LLaVA-Bench comparing our method’s response with the response generated by regular decoding using the LLaVA-1.5 model. Specifically, regular decoding often leads to hallucinated or inaccurate content, such as describing “taxi appears to be converted laundry machines” and “another person can be seen standing nearby”. In contrast, our response is more detailed, focusing on the fact that “a person is ironing clothes while on the move, which is an unconventional way”. The GPT-4V-aided evaluation in Table 4 further validates that our method improves both the accuracy and detailedness of generated responses.

## 4.3. Efficiency Comparison

In Table 5, we evaluate the efficiency of our approach using the LLaVA-1.5 model on the CHAIR benchmark, with a maximum token length of 128. We also report the performance of all compared methods across 5 benchmarks. Our approach demonstrates consistently superior performance, with only a  $1.07\times$  increase in time consumption and negligible additional GPU memory usage. These results validate that our approach is both efficient and effective, offering a

Method	Avg. Latency ↓	GPU Memory ↓	$\text{CHAIR}_S \downarrow$	MME ↑	POPE ↑	MMBench ↑	MM-Vet ↑
Regular	3.47 s ( $\times 1.00$ )	14945 MB ( $\times 1.00$ )	55.0	562.5	83.44	64.1	26.1
VCD	6.97 s ( $\times 2.01$ )	15749 MB ( $\times 1.05$ )	54.4	580.6	87.15	64.6	30.9
M3ID	7.05 s ( $\times 2.03$ )	15575 MB ( $\times 1.04$ )	54.4	598.1	87.52	64.4	29.9
OPERA	24.70 s ( $\times 7.12$ )	22706 MB ( $\times 1.52$ )	52.6	598.3	88.85	64.4	32.0
Woodpecker	10.68 s ( $\times 3.08$ )	22199 MB ( $\times 1.49$ )	57.6	588.3	86.45	64.0	30.6
HALC	22.61 s ( $\times 6.52$ )	23084 MB ( $\times 1.54$ )	51.0	579.6	87.68	64.2	30.8
<b>Ours</b>	3.70 s ( $\times 1.07$ )	14951 MB ( $\times 1.00$ )	<b>49.8</b>	<b>635.6</b>	<b>89.10</b>	<b>65.0</b>	<b>32.8</b>

Table 5. **Efficiency comparison.** For each method, we present the average inference latency per instance and peak GPU memory. Experiments are conducted on a single RTX A6000 Ada GPU.

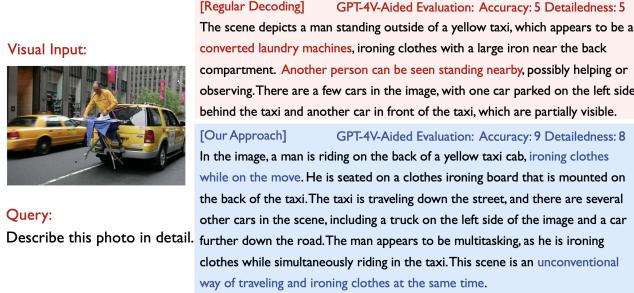


Figure 6. **Case study on the LLaVA-Bench benchmark.** We compare the responses generated by regular decoding and our method using LLaVA-1.5. GPT-4V-aided evaluation results are also provided alongside the responses. Hallucinated and accurate content is highlighted in red and blue.

favorable performance-cost trade-off.

#### 4.4. Ablation Study

**Selection of Layer for Textual Enhancement.** To investigate the impacts of choosing different layers for textual enhancement, we conduct ablation experiments on the POPE benchmark. Results in Figure 7 demonstrate that by selecting the initial layer for textual enhancement, our ONLY method achieves optimal performance on the POPE benchmark. Additionally, we observe that the performance of our approach is robust across different layers chosen for intervention, with ONLY exhibiting minimal variation and consistently outperforming VCD and M3ID. This robustness is due to our attention-head selection strategy, which dynamically selects different sets of heads across multiple layers, efficiently and effectively capturing language bias.

**Other Strategies for Textual Enhancement.** In Table 6, we compare the performance achieved by various textual enhancement strategies. Our approach of attention head selection using TVER achieves the best performance. In contrast, directly modifying attention weights—such as zeroing out or adding noise to visual attention weights, or doubling textual attention weights—results in suboptimal outcomes. Additionally, selecting attention heads based on the ratio of the sum of attention weights also leads to a performance decrease of 0.71% on POPE and 3.1% on CHAIR.

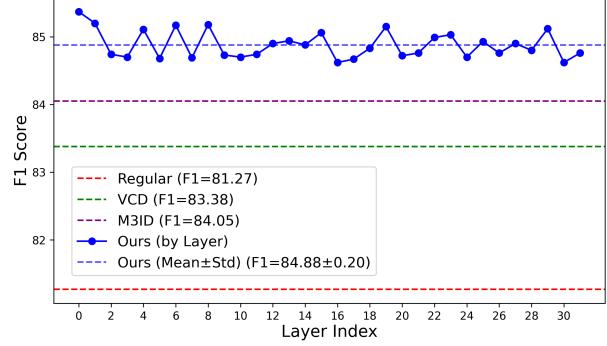


Figure 7. **Impacts of different selected layers.** We present the results obtained by selecting different layers for textual enhancement on the POPE benchmark using all 9,000 samples from COCO.

Strategy	POPE ↑				CHAIR ↓	
	Acc.	Prec.	Rec.	F1	$\text{CHAIR}_S$	$\text{CHAIR}_I$
Regular	80.42	78.20	84.59	81.27	26.2	9.4
$a_{t,i}^V \leftarrow 0$	84.26	82.13	87.69	84.82	21.2	6.9
$a_{t,i}^V \leftarrow a_{t,i}^V + \varepsilon$	83.95	81.67	88.16	84.79	22.1	7.6
$a_{t,i}^T \leftarrow a_{t,i}^T * 2$	84.37	82.52	87.55	84.96	21.6	6.8
Ratio $\leftarrow \sum a_T / \sum a_V$	84.20	81.57	87.56	84.46	23.1	8.2
<b>Ours</b>	<b>84.91</b>	<b>82.84</b>	<b>88.07</b>	<b>85.37</b>	<b>20.0</b>	<b>6.2</b>

Table 6. **Different Strategies for textual enhancement.** We conduct experiments with different textual enhancement strategies.

## 5. Conclusion

In this work, we introduce ONLY, a novel training-free approach that leverages a single additional Transformer layer to mitigate hallucinations in Large Vision-Language Models (LVLMs). By utilizing text-to-visual entropy, ONLY selectively activates attention heads with a high language bias to generate textually-enhanced outputs. These outputs are then adaptively decoded alongside the original output, using either contrastive or collaborative decoding. Extensive evaluations across six benchmarks and three LVLM backbones show that ONLY consistently outperforms existing methods in reducing hallucinations. Moreover, our approach incurs minimal additional inference time and memory consumption, making it well-suited for real-world applications that require real-time responses.

## Acknowledgements

This work has been funded in part by the Army Research Laboratory (ARL) award W911NF-23-2-0007 and W911QX-24-F-0049, DARPA award FA8750-23-2-1015, and ONR award N00014-23-1-2840. MM is supported by Meta AI Mentorship program and the research is partially supported by National Institutes of Health awards R01MH125740, R01MH132225, and R21MH130767.

## References

- [1] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- [2] Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou. Hallucination of multimodal large language models: A survey. *arXiv preprint arXiv:2404.18930*, 2024.
- [3] Jiawei Chen, Dingkang Yang, Tong Wu, Yue Jiang, Xiaolu Hou, Mingcheng Li, Shunli Wang, Dongling Xiao, Ke Li, and Lihua Zhang. Detecting and evaluating medical hallucinations in large vision language models. *arXiv preprint arXiv:2406.10185*, 2024.
- [4] Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. 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, 2024.
- [5] Zhiyang Chen, Yousong Zhu, Yufei Zhan, Zhaowen Li, Chaoyang Zhao, Jinqiao Wang, and Ming Tang. Mitigating hallucination in visual language models with visual supervision. *arXiv preprint arXiv:2311.16479*, 2023.
- [6] Zhaorun Chen, Zhuokai Zhao, Hongyin Luo, Huaxiu Yao, Bo Li, and Jiawei Zhou. HALC: Object hallucination reduction via adaptive focal-contrast decoding. In *International Conference on Machine Learning*, pages 7824–7846. PMLR, 2024.
- [7] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhang-hao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yong-hao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. 2023.
- [8] Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R. Glass, and Pengcheng He. Dola: Decoding by contrasting layers improves factuality in large language models. In *International Conference on Learning Representations*, 2024.
- [9] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: towards general-purpose vision-language models with instruction tuning. *Advances in Neural Information Processing Systems*, 36:49250–49267, 2023.
- [10] Alessandro Favero, Luca Zancato, Matthew Trager, Siddharth Choudhary, Pramuditha Perera, Alessandro Achille, Ashwin Swaminathan, and Stefano Soatto. Multi-modal hal-lucination control by visual information grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14303–14312, 2024.
- [11] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- [12] Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. Opera: Alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13418–13427, 2024.
- [13] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6700–6709, 2019.
- [14] Chaoya Jiang, Haiyang Xu, Mengfan Dong, Jiaxing Chen, Wei Ye, Ming Yan, Qinghao Ye, Ji Zhang, Fei Huang, and Shikun Zhang. Hallucination augmented contrastive learning for multimodal large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27036–27046, 2024.
- [15] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024.
- [16] Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. Mitigating object hal-lucinations 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, 2024.
- [17] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International Conference on Machine Learning*, pages 19730–19742. PMLR, 2023.
- [18] Xiang Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori B Hashimoto, Luke Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pages 12286–12312, 2023.
- [19] 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 Conference on Empirical Methods in Natural Language Processing*, pages 292–305, 2023.
- [20] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014.

- [21] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in Neural Information Processing Systems*, 36:34892–34916, 2023.
- [22] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024.
- [23] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024.
- [24] Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li, and Wei Peng. A survey on hallucination in large vision-language models. *arXiv preprint arXiv:2402.00253*, 2024.
- [25] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European Conference on Computer Vision*, pages 216–233. Springer, 2024.
- [26] Xinyu Lyu, Beita Chen, Lianli Gao, Hengtao Shen, and Jingkuan Song. Alleviating hallucinations in large vision-language models through hallucination-induced optimization. *Advances in Neural Information Processing Systems*, 37:122811–122832, 2024.
- [27] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [28] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 4035–4045, 2018.
- [29] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A benchmark for visual question answering using world knowledge. In *European Conference on Computer Vision*, pages 146–162. Springer, 2022.
- [30] Claude E Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948.
- [31] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023.
- [32] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9568–9578, 2024.
- [33] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [34] Zifu Wan, Yaqi Xie, Ce Zhang, Zhiqiu Lin, Zihan Wang, Simon Stepputtis, Deva Ramanan, and Katia Sycara. Instruction-part: Task-oriented part segmentation with instruction reasoning. *arXiv preprint arXiv:2505.18291*, 2025.
- [35] Xintong Wang, Jingheng Pan, Liang Ding, and Chris Biemann. Mitigating hallucinations in large vision-language models with instruction contrastive decoding. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 15840–15853, 2024.
- [36] Yixuan Wu, Yizhou Wang, Shixiang Tang, Wenhao Wu, Tong He, Wanli Ouyang, Philip Torr, and Jian Wu. Det-toolchain: A new prompting paradigm to unleash detection ability of mllm. In *European Conference on Computer Vision*, pages 164–182. Springer, 2024.
- [37] Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, and Fei Huang. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13040–13051, 2024.
- [38] 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. *arXiv preprint arXiv:2310.16045*, 2023.
- [39] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. MM-vet: Evaluating large multimodal models for integrated capabilities. In *International Conference on Machine Learning*, 2024.
- [40] Yuhang Zang, Wei Li, Jun Han, Kaiyang Zhou, and Chen Change Loy. Contextual object detection with multimodal large language models. *International Journal of Computer Vision*, pages 1–19, 2024.
- [41] Ce Zhang, Zifu Wan, Zhehan Kan, Martin Q Ma, Simon Stepputtis, Deva Ramanan, Russ Salakhutdinov, Louis-Philippe Morency, Katia P Sycara, and Yaqi Xie. Incorporating generative feedback for mitigating hallucinations in large vision-language models. In *Workshop on Responsibly Building the Next Generation of Multimodal Foundational Models*, 2024.
- [42] Ce Zhang, Kaixin Ma, Tianqing Fang, Wenhao Yu, Hongming Zhang, Zhisong Zhang, Yaqi Xie, Katia Sycara, Haitao Mi, and Dong Yu. Vscan: Rethinking visual token reduction for efficient large vision-language models. *arXiv preprint arXiv:2505.22654*, 2025.
- [43] Ce Zhang, Zifu Wan, Zhehan Kan, Martin Q. Ma, Simon Stepputtis, Deva Ramanan, Russ Salakhutdinov, Louis-Philippe Morency, Katia P. Sycara, and Yaqi Xie. Self-correcting decoding with generative feedback for mitigating hallucinations in large vision-language models. In *International Conference on Learning Representations*, 2025.
- [44] Jinrui Zhang, Teng Wang, Haigang Zhang, Ping Lu, and Feng Zheng. Reflective instruction tuning: Mitigating hallucinations in large vision-language models. In *European Conference on Computer Vision*, 2024.

# ONLY: One-Layer Intervention Sufficiently Mitigates Hallucinations in Large Vision-Language Models

## Supplementary Material

This supplementary document is organized as follows:

- The intuitive and theoretical explanation for our motivation is provided in Section A.
- Additional experimental details, including further implementation details, descriptions of other implemented baselines, and license information for the utilized code and datasets, are provided in Section B.
- Additional experimental results on different benchmarks are presented in Section C.
- Additional ablation studies with different parameters are presented in Section D.
- More case studies and GPT-4V-aided evaluations are provided in Section E.
- Potential directions for future work are discussed in Section F.

### A. More Explanation on Motivation

#### A.1. Intuitive Explanation for TVER

Our method is motivated by a principle in information theory:  $H(x|y) \leq H(x)$ . Let  $H(\mathcal{T})$  and  $H(\mathcal{V})$  denote the entropy of pure textual and visual attention, respectively. During LVLM decoding, since the model processes both image and text simultaneously, we treat the attention distributions as conditioned on the other modality. This leads to an approximate theoretical form of Eq. (11):  $TVER = \frac{H(\mathcal{T}|\mathcal{V})}{H(\mathcal{V}|\mathcal{T})}$ . Since  $H(\mathcal{T}|\mathcal{V}) \leq H(\mathcal{T})$  and  $H(\mathcal{V}|\mathcal{T}) \leq H(\mathcal{V})$ , a higher  $H(\mathcal{T}|\mathcal{V})$  indicates behavior closer to purely textual inference, while higher  $H(\mathcal{V}|\mathcal{T})$  suggests reliance on visual priors. To approximate the noisy branch used in VCD and M3ID, we aim to enhance textual focus and suppress visual focus, which motivates maximizing TVER for effective textual enhancement.

### B. More Experimental Details

#### B.1. Benchmarks and Metrics

We conduct extensive experiments on the following benchmarks:

- **POPE** [19] is a popular benchmark for assessing object hallucinations in LVLMs. It tests the models with yes-or-no questions regarding the presence of specific objects, such as, “Is there a {object} in the image?” The images from the benchmark derive from three existing datasets: MSCOCO [20], A-OKVQA [29],

and GQA [13], and comprises three distinct subsets—*random*, *popular*, and *adversarial*—based on how the negative samples are generated. For each dataset setting, the benchmark provides 6 questions per image, resulting in 3,000 test instances. We evaluate the performance of different methods using four metrics: accuracy, precision, recall, and F1 score.

- **CHAIR** [28] evaluates object hallucinations through image captioning, where the LVLMs are prompted to describe 500 randomly selected images from the MSCOCO validation set. The performance is evaluated based on two metrics:

$$\text{CHAIR}_I = \frac{\# \text{ hallucinated objects}}{\# \text{ all objects mentioned}}, \quad (21)$$

$$\text{CHAIR}_S = \frac{\# \text{ sentences with hallucinated object}}{\# \text{ all sentences}}. \quad (22)$$

- **MME-Hallucination** [11] is a comprehensive benchmark consisting of four subsets: *existence* and *count* for object-level hallucinations, and *position* and *color* for attribute-level hallucinations. Each subset includes 30 images and 60 questions, with two questions per image. Similar to POPE [19], the benchmark includes yes-or-no questions, and performance is assessed based on binary accuracy. Following the official implementation, the reported score is calculated by combining accuracy and accuracy+, where accuracy is based on individual questions, and accuracy+ is based on images where both questions are answered correctly.

- **MMBench** [25] is a comprehensive benchmark designed to evaluate LVLMs’ multimodal understanding and reasoning abilities. It emphasizes tasks that require integrating visual and textual information, assessing a model’s performance in diverse, real-world scenarios. MMBench employs a hierarchical ability taxonomy, categorizing Perception and Reasoning as Level-1 (L-1) abilities. This taxonomy is further refined into six Level-2 (L-2) dimensions and twenty Level-3 (L-3) dimensions, providing a detailed framework for assessment.

- **MMVP** [32] is a benchmark designed to assess the fine-grained visual recognition capabilities of LVLMs using CLIP-blind pairs. It comprises 150 image pairs, each paired with a binary-option question. Each image is evaluated separately, and an LVLM’s response is deemed correct only if it answers both questions associated with a pair accurately.

- **MM-Vet** [39] is a benchmark for evaluating LVLMs on complex tasks. It defines 6 core vision-language capabil-

ties, including recognition, OCR, knowledge, language generation, spatial awareness, and math. An LLM-based evaluator is used to ensure consistent evaluation across diverse question types. The dataset includes 187 images from various online sources and collects 205 questions, each of which requires one or more capabilities to answer.

- **LLaVA-Bench**<sup>1</sup> includes 24 images depicting complex scenes, memes, paintings, and sketches, accompanied by 60 challenging questions. Selected examples from this dataset are used for qualitative comparisons of responses generated by different decoding methods. Additionally, following Yin et al. [38], we evaluate the accuracy and level of detail in the generated responses using the advanced LVLM, GPT-4V<sup>2</sup>.

## B.2. More Implementation Details

In our experiments, we adhere to the default query format for the input data used in both LLaVA-1.5 [21], InstructBLIP [9], and Qwen-VL [1]. We set  $\alpha_1 = 3$ ,  $\alpha_2 = 1$  by default in our decoding process. Additionally, we set  $\gamma = 0.2$  for LLaVA-1.5 and  $\gamma = 0.4$  for InstructBLIP/Qwen-VL. We follow VCD [16] to implement adaptive plausibility constraints [18]:

$$p_\theta(y_t) = 0, \quad \text{if } y_t \notin \mathcal{S}(y_{<t}), \quad (23)$$

where  $\mathcal{S}(y_{<t}) = \{y_t \in \mathcal{S} : p_\theta(y_t|v, \mathbf{x}, \mathbf{y}_{<t}) \geq \beta \max_w p_\theta(w|v, \mathbf{x}, \mathbf{y}_{<t})\}$ . Here,  $\mathcal{S}$  is the whole vocabulary of LVLM, and hyperparameter  $\beta \in [0, 1]$  controls the truncation of the next token distribution. A larger  $\beta$  indicates more aggressive truncation, keeping only the high-probability tokens. In our implementation, we set the logits for  $y_t \notin \mathcal{S}(y_{<t})$  to  $-\infty$ . By default, we set  $\beta = 0.1$  for all tasks. All experiments are conducted on a single 48GB NVIDIA RTX 6000 Ada GPU.

## B.3. Pilot Study Details

For Figure 4, we visualize 500 images from the CHAIR [28] benchmark (left) and 3,000 images from POPE [19] (right). For Figure 3, we analyze 3,000 POPE images to examine the relationship between entropy deviation and noise level.

## B.4. Devision of Textual and Visual Tokens

In Eq. 8, textual and visual attention are obtained based on the indices corresponding to each modality. The index ranges for both modalities are listed below:

- LLaVA-1.5 [21]:  
Textual indices – [0:35], [611:]; Visual indices – [35:611].
- InstructBLIP [9]:  
Textual indices – [32:]; Visual indices – [0:32].

<sup>1</sup><https://huggingface.co/datasets/liuhaojian/llava-bench-in-the-wild>.

<sup>2</sup><https://openai.com/index/gpt-4v-system-card>.

- Qwen-VL[1]:  
Textual indices – [257:]; Visual indices – [1:257].

## B.5. Details of Other Baselines

In this work, we mainly compare the performance of our ONLY with two state-of-the-art contrastive-decoding approaches: VCD [16] and M3ID [10]. The method and implementation details for these approaches are provided below:

- **VCD** [16] contrasts output distributions derived from original and distorted visual inputs. Specifically, given a textual query  $x$  and a visual input  $v$ , the model generates two distinct output distributions: one conditioned on the original  $v$  and the other conditioned on the distorted visual input  $v'$ , which is obtained by applying pre-defined distortions (e.g., Gaussian noise mask) to  $v$ . Then, a new contrastive probability distribution is computed as:

$$p_{vcd}(y_t) = \text{softmax}[(1 + \alpha)f_\theta(y|v, \mathbf{x}, \mathbf{y}_{<t}) - \alpha f_\theta(y|v', \mathbf{x}, \mathbf{y}_{<t})]. \quad (24)$$

In our implementation, we follow the default setting in VCD [16] and set  $\alpha = 1$  for reproduction. To generate  $v'$ , we use a total of 500 noise steps.

- **M3ID** [10] contrasts output distributions derived from original visual inputs with those from pure text inputs, which lack visual information. The final probability distribution is given by:

$$p_{m3id}(y_t) = \text{softmax}[f_\theta(y|v, \mathbf{x}, \mathbf{y}_{<t}) + \frac{1 - e^{-\lambda t}}{e^{-\lambda t}} (f_\theta(y|v, \mathbf{x}, \mathbf{y}_{<t}) - f_\theta(y|\mathbf{x}, \mathbf{y}_{<t}))]. \quad (25)$$

Following their recommended best practice, we set the hyperparameter  $\lambda$ , which balances the conditioned and unconditioned models, to 0.02.

## B.6. Dataset and Code Licensing

**Datasets.** We list the known license information for the datasets below: POPE [19] and MMVP [32] benchmarks are licensed under MIT License. CHAIR [28] is made available under the BSD 2-Clause License. LLaVA-Bench is available under Apache-2.0 License. MME-Hallucination [11] benchmark dataset is collected by Xiamen University for academic research only. MM-Vet [39] dataset is under the CC BY-NC 4.0 license.

**Code.** In this work, we also use some code implementations from the existing codebases: LLaVA [21] and VCD [16] are licensed under the Apache-2.0 License. InstructBLIP [9] is under BSD-3-Clause License. Qwen-VL [1] is under the Tongyi Qianwen License.

Table C1. **Results on MME-Hallucination [11] benchmark.** We report the average MME scores along with the standard deviation across three random seeds for each subset. We also report the total scores achieved by the different methods across all four subsets in the final column. Higher scores ( $\uparrow$ ) indicate better performance. The best results are **bolded**, and the second-best are underlined.

Model	Method	Object-level		Attribute-level		Total Score $\uparrow$
		Existence $\uparrow$	Count $\uparrow$	Position $\uparrow$	Color $\uparrow$	
LLaVA-1.5	Regular	173.75 ( $\pm 4.79$ )	121.67 ( $\pm 12.47$ )	117.92 ( $\pm 3.69$ )	149.17 ( $\pm 7.51$ )	562.50 ( $\pm 3.96$ )
	DoLa	176.67 ( $\pm 2.89$ )	113.33 ( $\pm 10.41$ )	90.55 ( $\pm 8.22$ )	141.67 ( $\pm 7.64$ )	522.22 ( $\pm 16.78$ )
	OPERA	183.33 ( $\pm 6.45$ )	<u>137.22</u> ( $\pm 6.31$ )	122.78 ( $\pm 2.55$ )	<u>155.00</u> ( $\pm 5.00$ )	<u>598.33</u> ( $\pm 10.41$ )
	VCD	186.67 ( $\pm 5.77$ )	125.56 ( $\pm 3.47$ )	128.89 ( $\pm 6.73$ )	139.45 ( $\pm 12.51$ )	580.56 ( $\pm 15.13$ )
	M3ID	186.67 ( $\pm 5.77$ )	128.33 ( $\pm 10.41$ )	<u>131.67</u> ( $\pm 5.00$ )	151.67 ( $\pm 20.88$ )	598.11 ( $\pm 20.35$ )
	Woodpecker	<u>187.50</u> ( $\pm 2.89$ )	125.00 ( $\pm 0.00$ )	126.66 ( $\pm 2.89$ )	149.17 ( $\pm 17.34$ )	588.33 ( $\pm 10.00$ )
	HALC	183.33 ( $\pm 0.00$ )	133.33 ( $\pm 5.77$ )	107.92 ( $\pm 3.69$ )	<u>155.00</u> ( $\pm 5.00$ )	579.58 ( $\pm 9.07$ )
InstructBLIP	<b>Ours</b>	<b>191.67</b> ( $\pm 2.89$ )	<b>145.55</b> ( $\pm 10.72$ )	<b>136.66</b> ( $\pm 2.89$ )	<b>161.66</b> ( $\pm 2.89$ )	<b>635.55</b> ( $\pm 5.85$ )
Qwen-VL	Regular	160.42 ( $\pm 5.16$ )	79.17 ( $\pm 8.22$ )	<b>79.58</b> ( $\pm 8.54$ )	<u>130.42</u> ( $\pm 17.34$ )	449.58 ( $\pm 24.09$ )
	DoLa	175.00 ( $\pm 5.00$ )	55.00 ( $\pm 5.00$ )	48.89 ( $\pm 3.47$ )	113.33 ( $\pm 6.67$ )	392.22 ( $\pm 7.88$ )
	OPERA	175.00 ( $\pm 3.33$ )	61.11 ( $\pm 3.47$ )	53.89 ( $\pm 1.92$ )	120.55 ( $\pm 2.55$ )	410.56 ( $\pm 9.07$ )
	VCD	158.89 ( $\pm 5.85$ )	<b>91.67</b> ( $\pm 18.34$ )	66.11 ( $\pm 9.76$ )	121.67 ( $\pm 12.58$ )	438.33 ( $\pm 16.07$ )
	M3ID	160.00 ( $\pm 5.00$ )	<u>87.22</u> ( $\pm 22.63$ )	69.44 ( $\pm 9.18$ )	125.00 ( $\pm 7.64$ )	441.67 ( $\pm 17.32$ )
	<b>Ours</b>	<b>180.00</b> ( $\pm 5.00$ )	77.78 ( $\pm 7.70$ )	<u>74.44</u> ( $\pm 12.05$ )	<b>135.55</b> ( $\pm 3.85$ )	<b>467.77</b> ( $\pm 8.55$ )
Qwen-VL	Regular	155.00 ( $\pm 3.54$ )	127.67 ( $\pm 13.36$ )	131.67 ( $\pm 7.73$ )	173.00 ( $\pm 9.75$ )	587.33 ( $\pm 31.06$ )
	VCD	156.00 ( $\pm 6.52$ )	131.00 ( $\pm 6.19$ )	128.00 ( $\pm 3.61$ )	<b>181.67</b> ( $\pm 5.14$ )	596.67 ( $\pm 11.61$ )
	M3ID	<u>178.33</u> ( $\pm 2.89$ )	<u>143.33</u> ( $\pm 2.89$ )	<u>150.00</u> ( $\pm 2.89$ )	175.00 ( $\pm 5.00$ )	<u>646.66</u> ( $\pm 8.50$ )
	<b>Ours</b>	<b>180.00</b> ( $\pm 5.00$ )	<b>146.67</b> ( $\pm 5.00$ )	<b>156.11</b> ( $\pm 6.31$ )	<u>178.33</u> ( $\pm 2.89$ )	<b>661.11</b> ( $\pm 3.47$ )

Method	LR	AR	RR	FP-S	FP-C	CP	Overall
Regular	30.51	71.36	52.17	67.58	<b>58.74</b>	76.35	64.09
VCD	30.51	<b>73.37</b>	53.04	<b>67.92</b>	57.34	77.03	64.60
M3ID	30.51	72.36	53.04	67.58	57.34	<b>77.36</b>	64.43
<b>Ours</b>	<b>33.05</b>	<b>73.37</b>	<b>54.78</b>	66.55	<b>58.74</b>	<b>77.36</b>	<b>64.95</b>

Table C2. **Detailed results on MMBench benchmark.** Abbreviations adopted: LR for Logical Reasoning; AR for Attribute Reasoning; RR for Relation Reasoning; FP-S for Fine-grained Perception (Single Instance); FP-C for Fine-grained Perception (Cross Instance); CP for Coarse Perception. The best results are **bolded**.

Method	Rec	OCR	Know	Gen	Spat	Math	Total
Regular	30.8	19.0	14.5	17.9	26.9	<b>11.5</b>	26.1
VCD	35.6	21.9	18.3	<u>21.9</u>	28.9	3.8	30.9
M3ID	35.0	19.7	18.8	<u>19.0</u>	26.0	7.7	29.9
DoLA	<u>37.2</u>	22.1	17.9	21.0	26.3	7.7	31.7
OPERA	35.4	<b>25.6</b>	<u>20.5</u>	<b>22.9</b>	30.9	11.5	<u>32.0</u>
HALC	36.2	21.5	17.5	20.1	23.5	<b>7.7</b>	30.8
<b>Ours</b>	<b>37.3</b>	<u>23.9</u>	<b>22.9</b>	<u>22.1</u>	<b>31.3</b>	3.8	<b>32.8</b>

Table C3. **Detailed results on MM-Vet benchmark.** Abbreviations adopted: Rec for Recognition, OCR for Optical Character Recognition, Know for Knowledge, Gen for Language Generation, Spat for Spatial Awareness, Math for Mathematics. The best results are **bolded**, and the second best are underlined.

## C. More Experimental Results and Analysis

### C.1. Full Results on MME-Hallucination

In Table C1, we present the full results on the MME-Hallucination benchmark. From the results, our method consistently outperforms others on both object-level and attribute-level data across three LVLM backbones.

### C.2. Full Results on MMBench

In Table C2, we present the overall performance on the MMBench benchmark, as well as the detailed performance across six Level-2 abilities: Logical Reasoning (LR), Attribute Reasoning (AR), Relation Reasoning (RR), Fine-grained Perception - Single Instance (FP-S), Fine-grained

Perception - Cross Instance (FP-C), and Coarse Perception (CP). We follow VCD [16] to conduct experiments on the MMBench-dev set. Our method outperforms other baselines in most abilities and the overall score.

### C.3. Results on MM-Vet

In Table C3, we present the overall performance on the MM-Vet [39] benchmark, where we use LLaVA-1.5 as the LVLM backbone. From the results, we observed that our method consistently outperforms others on the MM-Vet benchmark.

#### C.4. Evaluation on other advanced LVLMs

We further report results of LLaVA-NeXT-7B/13B [23] on POPE (MS-COCO) benchmark in table C4. Our method consistently outperforms existing approaches at both scales while requiring only half the inference time and resources.

Method	LLaVA-NeXT-7B				LLaVA-NeXT-13B			
	Acc.	Prec.	Rec.	F1	Acc.	Prec.	Rec.	F1
Regular	85.71	85.27	86.33	85.80	86.74	86.53	87.04	86.78
VCD	87.07	87.40	86.62	87.01	87.09	87.39	86.69	87.04
M3ID	87.48	87.64	<b>87.27</b>	87.45	87.84	<b>87.95</b>	87.71	87.83
Ours	<b>87.96</b>	<b>88.59</b>	87.13	<b>87.86</b>	<b>87.94</b>	87.31	<b>88.80</b>	<b>88.05</b>

Table C4. **Detailed results with LLaVA-NeXT**. The best results are **bolded**, and the second best are underlined.

## D. More Ablation Studies and Analysis

### D.1. Effects of $\alpha_1$ and $\alpha_2$ in Adaptive Decoding

In Section 3, we introduce collaborative and contrastive decoding, along with hyperparameters  $\alpha_1$  and  $\alpha_2$ , which regulate the influence of the textual-enhanced branch. Tables D5 and D6 analyze their impact, showing that the default values  $\alpha_1 = 3$  and  $\alpha_2 = 1$  yield the best performance across benchmarks. Notably, setting these to 0 reduces our approach to standard decoding, confirming that adaptive decoding significantly enhances hallucination mitigation in LVLMs.

### D.2. Effect of $\beta$ in Adaptive Plausibility Constraint

We perform an ablation study on  $\beta$ , introduced in Eq. 23, by varying its value from 0 to 0.5 while keeping all other hyperparameters fixed. As shown in Table D7, setting  $\beta = 0$ , which removes the constraint, leads to suboptimal performance across both benchmarks. Our method achieves the best results with  $\beta = 0.1$ , which we adopt as the default setting.

### D.3. Effect of $\gamma$ in Adaptive Plausibility Constraint

We further studied the influence led by the threshold  $\gamma$  for adaptive decoding. The results in Table D8 show that setting  $\gamma = 0.2$  reaches the optimal result for LLaVA-1.5. Besides, we keep  $\gamma = 0.4$  for other baseline LVLMs.

### D.4. Scaling Up the LVLMs

We extend our evaluation to the 13B variant of the LLaVA-1.5 model to assess the scalability of our approach. Table D9 compares our results with state-of-the-art methods across all three subsets of the POPE benchmark using the 13B model. Our findings show that increasing model size does not mitigate hallucination issues, as the 7B and 13B models exhibit comparable performance. Notably, ONLY consistently outperforms other approaches across all subsets, demonstrating its effectiveness and scalability.

Values	POPE				CHAIR	
	Acc.	Prec.	Rec.	F1	CHAIR <sub>S</sub>	CHAIR <sub>I</sub>
$\alpha_1 = 0$	88.13	<b>94.55</b>	80.93	87.21	23.5	8.6
$\alpha_1 = 1$	88.27	94.50	81.27	87.38	22.4	7.8
$\alpha_1 = 2$	88.87	89.63	88.10	88.86	21.5	7.2
$\alpha_1 = 3$	<b>89.70</b>	89.95	<b>88.27</b>	<b>89.10</b>	<b>20.0</b>	<b>6.2</b>
$\alpha_1 = 4$	88.37	88.85	87.94	88.39	22.3	7.6

Table D5. **Sensitivity analysis of hyperparameter  $\alpha_1$** . We present the performance of our approach, based on the LLaVA-1.5 backbone, across two benchmarks for varying values of  $\alpha_1$ . Note that we fix  $\alpha_2 = 1$  in this experiment.

Values	POPE				CHAIR	
	Acc.	Prec.	Rec.	F1	CHAIR <sub>S</sub>	CHAIR <sub>I</sub>
$\alpha_2 = 0$	86.50	86.35	88.13	86.72	24.8	9.3
$\alpha_2 = 1$	<b>89.70</b>	89.95	<b>88.27</b>	<b>89.10</b>	<b>20.0</b>	<b>6.2</b>
$\alpha_2 = 2$	87.67	96.69	78.00	86.35	22.4	7.6
$\alpha_2 = 3$	87.37	<b>97.14</b>	77.00	85.91	23.4	7.3
$\alpha_2 = 4$	87.13	97.12	76.53	85.61	24.2	8.1

Table D6. **Sensitivity analysis of hyperparameter  $\alpha_2$** . We present the performance of our approach, based on the LLaVA-1.5 backbone, across two benchmarks for varying values of  $\alpha_2$ . Note that we fix  $\alpha_1 = 3$  in this experiment.

Values	POPE				CHAIR	
	Acc.	Prec.	Rec.	F1	CHAIR <sub>S</sub>	CHAIR <sub>I</sub>
$\beta = 0$	87.70	93.40	81.13	86.84	24.6	10.1
$\beta = 0.05$	88.17	<b>94.21</b>	81.33	87.30	23.7	9.6
$\beta = 0.1$	<b>89.70</b>	89.95	<b>88.27</b>	<b>89.10</b>	<b>20.0</b>	<b>6.2</b>
$\beta = 0.25$	89.56	89.48	87.63	88.55	21.4	7.6
$\beta = 0.5$	89.47	89.83	86.53	88.15	22.1	7.2

Table D7. **Sensitivity analysis of hyperparameter  $\beta$** . We present the performance of our approach, based on the LLaVA-1.5 backbone, across two benchmarks for varying values of  $\beta$ .

Values	POPE				CHAIR	
	Acc.	Prec.	Rec.	F1	CHAIR <sub>S</sub>	CHAIR <sub>I</sub>
$\gamma = 0.0$	89.13	90.41	86.38	88.35	23.5	8.2
$\gamma = 0.1$	89.20	89.88	86.73	88.28	22.6	8.1
$\gamma = 0.2$	<b>89.70</b>	89.95	<b>88.27</b>	<b>89.10</b>	<b>20.0</b>	<b>6.2</b>
$\gamma = 0.3$	89.40	93.20	85.00	88.91	21.2	7.1
$\gamma = 0.4$	89.03	93.99	83.40	88.38	21.7	7.0
$\gamma = 0.5$	89.15	92.26	84.29	88.10	22.4	7.6
$\gamma = 0.6$	89.21	91.78	85.39	88.47	23.1	8.1

Table D8. **Sensitivity analysis of hyperparameter  $\gamma$** . We present the performance of our approach, based on the LLaVA-1.5 backbone, across two benchmarks for varying values of  $\gamma$ .

### D.5. Details about Ablation Studies on Layer Selection and Strategies

In Section 4.4, we conduct two ablation studies to validate our proposed method. Detailed results are provided

Table D9. **Results on POPE [19] benchmark using 13B-sized LLaVA-1.5.** Higher ( $\uparrow$ ) accuracy, precision, recall, and F1 indicate better performance.

Setup	Method	LLaVA-1.5			
		Acc. $\uparrow$	Prec. $\uparrow$	Rec. $\uparrow$	F1 $\uparrow$
Random	Regular	82.53	78.57	89.47	83.67
	VCD	84.80	80.67	91.53	85.76
	M3ID	85.37	81.30	91.87	86.26
	<b>Ours</b>	<b>88.63</b>	<b>89.66</b>	87.33	<b>88.48</b>
Popular	Regular	80.53	76.17	88.87	82.03
	VCD	82.23	76.88	92.20	83.84
	M3ID	82.60	77.91	91.00	83.95
	<b>Ours</b>	<b>85.47</b>	<b>83.25</b>	88.80	<b>85.94</b>
Adversarial	Regular	75.80	70.41	89.00	78.62
	VCD	77.33	71.44	91.07	80.07
	M3ID	77.43	71.65	90.80	80.09
	<b>Ours</b>	<b>80.63</b>	<b>76.33</b>	88.80	<b>82.10</b>

below.

**Selection of Layer for Textual Enhancement:** In this experiment, we select a single layer from the total of 32 layers for textual enhancement. The F1 scores for our method across the 32 layers are as follows: [85.37, 85.20, 84.74, 84.7, 85.11, 84.68, 85.17, 84.69, 85.18, 84.73, 84.7, 84.74, 84.9, 84.94, 84.88, 85.06, 84.62, 84.67, 84.83, 85.15, 84.72, 84.76, 84.99, 85.03, 84.7, 84.93, 84.76, 84.9, 84.8, 85.12, 84.62, 84.76]. In comparison, the results for regular decoding, VCD [16], and M3ID [10] are 81.27, 83.38, and 84.05, respectively.

**Other Strategies for Textual Enhancement:** In Table 6, we explore additional strategies for textual enhancement, which include:

- $a_{\ell,i}^V \leftarrow 0$ : Setting the visual attention in the attention matrix to zero, inspired by M3ID [10], which uses a visual-free input for contrastive decoding;
- $a_{\ell,i}^V \leftarrow a_{\ell,i}^V + \varepsilon$ : Adding noise  $\varepsilon$  to the visual attention, inspired by VCD [16], which uses a distorted visual input for contrastive decoding;
- $a_{\ell,i}^T \leftarrow a_{\ell,i}^T * 2$ : Enhancing textual attention by directly multiplying it by 2;
- Ratio  $\leftarrow \sum a_T / \sum a_V$ : Instead of using the text-to-visual entropy ratio as the criterion to select textual-enhanced heads, we use the ratio between the sum of textual attention and visual attention. Heads with a ratio lower than the average across all heads are masked out, as described in Eq. 12.

All of these strategies require minimal additional computation, providing an efficiency advantage over other methods [10, 16]. This demonstrates the effectiveness of using just one layer for mitigating hallucinations in LVLMs, rather than relying on an extra full-process inference.

## E. More Case Studies

### E.1. Details about GPT-4V-Aided Evaluation

Following VCD [16], we use GPT-4V to evaluate responses in open-ended generation scenarios, scoring them based on accuracy and detailedness. Leveraging GPT-4V’s strong human-like capabilities, it can detect incorrect colors, positions, and relationships, allowing for a thorough evaluation of the responses.

Specifically, we apply the prompt in Table E10 to instruct GPT-4V to rate two responses on a scale from 1 to 10 for both accuracy and detailedness:

- **Accuracy** measures the consistency between the responses/descriptions generated by the LVLMs and the given image. A lower score is given if GPT-4V detects any inconsistencies in the content.
- **Detailedness** evaluates the depth and specificity of the responses. A higher score is awarded if the response includes comprehensive descriptions, captures fine-grained details of the image, and provides well-elaborated explanations. Conversely, a lower score is given if the response is vague or lacks sufficient detail.

### E.2. More Qualitative Results

In Figure E1, we present additional case studies from LLaVA-Bench to qualitatively demonstrate the effectiveness of our methods in mitigating hallucinations. We also include GPT-4V evaluations, providing accuracy and detailedness scores for each instance.

## F. Future Work

In future work, we aim to further improve the speed of our method and develop a more efficient hallucination mitigation approach that surpasses the original LVLM speed, leveraging efficient LVLM techniques like FastV [4] and VScan [42]. Additionally, we plan to explore our method’s potential for video hallucination mitigation to demonstrate its adaptability across various tasks.

---

**Description:**

AI that scores image description accuracy and detailedness.

---

**Instructions:**

You are an AI designed to evaluate and score the performance of two AI assistants in describing a given image. Your primary focus is on the accuracy and detailedness of their descriptions. You will assess the accuracy by checking for hallucinations - any part of the description that is inconsistent with the image content. For detailedness, you will consider how rich the response is in necessary details, excluding any hallucinated parts. You will provide scores on a scale from 1 to 10 for each assistant separately, based on these criteria. After scoring, you will offer an explanation for your evaluation, ensuring it is free from bias and not influenced by the order of presentation of the responses.

Input format:

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

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

Output format:

Accuracy:

Scores of the two answers:

Reason:

Detailedness:

Scores of the two answers:

Reason:

---

Table E10. **GPT-4V-aided evaluation setup.** We present the prompt we provided to GPT-4V to evaluate the LVLM responses based on accuracy and detailedness.

<p><b>Visual Input:</b></p> 	<p><b>[Regular Decoding]</b> GPT-4V-Aided Evaluation: Accuracy: 2 Detailedness: 3 The image features a myJokeWebsite box placed on a piece of graph paper. This graph paper is filled with multiple sentences, likely specifically written jokes. In addition to the main box, there are smaller boxes placed around the main page, primarily in the lower half of the paper. These smaller boxes contain various jokes, giving the impression of a collection of joke examples to choose from for each joke box.</p> <p><b>[Our Approach]</b> GPT-4V-Aided Evaluation: Accuracy: 4 Detailedness: 5 The image is a blue drawing of a joke website sign that says "My Joke Website: Funny Joke." Below the sign, there is a box labeled "Push to Reveal Punchline". The drawing appears to be done in blue ink on a piece of paper, representing the user's intention to access funny jokes on the website.</p>
<p><b>Query:</b> Describe this photo in detail.</p>	
<p><b>Visual Input:</b></p> 	<p><b>[Regular Decoding]</b> GPT-4V-Aided Evaluation: Accuracy: 5 Detailedness: 4 The meme features a cartoon character representing a person, pointing at a graph or animating figures on a large chart, possibly a white board, while surrounded by a dose of humor. The character appears to be reacting to an equation, discussing or teaching the presence. Below the main figure, four additional figures are portrayed, each interacting or holding their hands up in response to the main character's commentary about theory graphs. On the side, there is a description providing additional context.</p> <p><b>[Our Approach]</b> GPT-4V-Aided Evaluation: Accuracy: 7 Detailedness: 6 The image consists of two comics, one above the other, each featuring a person in a similar setting. In both comics, the person is standing in front of a chart or graph, which is likely discussing layers. The main focus of the comics is to emphasize the importance of layers and their applications in various fields. In each comic, there is a character standing in front of the chart, while another character is drawing attention to the chart. The drawings are simple, yet they convey the message effectively.</p>
<p><b>Query:</b> Describe this photo in detail.</p>	
<p><b>Visual Input:</b></p> 	<p><b>[Regular Decoding]</b> GPT-4V-Aided Evaluation: Accuracy: 4 Detailedness: 5 The image showcases a large painting of a monkey dog with a sword through its throat, surrounded by a Renaissance-style painting. In the scene, a woman resembling the famous Mona Lisa has a dog wearing a hood, and instead of her looking out of the canvas, it appears as if the dog is staring back. The dog is the focal point of the scene, adding a mysterious twist to a traditional painting.</p> <p><b>[Our Approach]</b> GPT-4V-Aided Evaluation: Accuracy: 8 Detailedness: 9 The painting depicts a unique scene where a dog is portrayed as a woman, possibly Mona Lisa, wearing a hood. The dog is adorned with a Renaissance-style dress and a hat, giving it an intriguing resemblance to a human. The background of the painting features mountains and water, further enhancing the artistic portrayal. A rock formation can also be seen in the painting. The overall composition creates a captivating and amusing artistic representation of the dog.</p>
<p><b>Query:</b> Describe this photo in detail.</p>	

**Figure E1. Case studies on the LLaVA-Bench benchmark.** We compare the responses generated by regular decoding and our method using LLaVA-1.5. GPT-4V-aided evaluation results are also provided alongside the responses. Hallucinated and accurate content is highlighted in red and blue.