# Multi-Frequency Contrastive Decoding: Alleviating Hallucinations for Large Vision-Language Models

**Bingqian Liu, Fu Zhang[*], Guoqing Chen, Jingwei Cheng**

School of Computer Science and Engineering, Northeastern University, China

{liubingqian.neu, chenguoqing247}@gmail.com
{zhangfu, chengjingwei}@neu.edu.cn

## Abstract

Large visual-language models (LVLMs) have demonstrated remarkable performance in visual-language tasks. However, object hallucination remains a significant challenge for LVLMs. Existing studies attribute object hallucinations in LVLMs mainly to linguistic priors and data biases. We further explore the causes of object hallucinations from the perspective of frequency domain and reveal that insufficient frequency information in images amplifies these linguistic priors, increasing the likelihood of hallucinations. To mitigate this issue, we propose the Multi-Frequency Contrastive Decoding (MFCD) method, a simple yet training-free approach that removes the hallucination distribution in the original output distribution, which arises from LVLMs neglecting the high-frequency information or low-frequency information in the image input. Without compromising the general capabilities of LVLMs, the proposed MFCD effectively mitigates the object hallucinations in LVLMs. Our experiments demonstrate that MFCD significantly mitigates object hallucination across diverse large-scale vision-language models, without requiring additional training or external tools. In addition, MFCD can be applied to various LVLMs without modifying model architecture or requiring additional training, demonstrating its generality and robustness. Codes are available at https://github.com/liubq-dev/mfcd.

## 1 Introduction

In recent years, large language models (LLMs) (Brown et al., 2020; Touvron et al., 2023; Yenduri et al., 2023; Xu et al., 2025; Qwen et al., 2025) have demonstrated outstanding performance in a variety of text-based tasks. Many researchers have been dedicated to extending the powerful language capabilities of LLMs to the visual domain. They have achieved this by combining the visual encoders

---

[*] Corresponding author.



Figure 1: An illustration of absence of frequency amplifying language priors. On the left is the input of the original image, and on the right is the input of the image removed the high-frequency information (i.e., the image information in the regions where the pixel values change drastically). The probabilities of generating a radio and a cup are respectively shown on the lower part. When images with high-frequency information removed are used as input, the probability of the hallucinated object "cup" as the next token is higher. MFCD mitigates linguistic bias in LVLMs' outputs by dynamically comparing these probabilities.

from CLIP (Radford et al., 2021) with LLMs to form LVLMs (Liu et al., 2023; Li et al., 2023a; Wang et al., 2024a; Li et al., 2024; Chen et al., 2025b; Bai et al., 2025), thereby extending the language capabilities of LLMs to the visual field. LVLMs perform excellently in various visual tasks, as well as in more complex tasks such as content understanding and generation (Geng et al., 2024).

However, LVLMs still face the challenge of object hallucination (Li et al., 2023c; Gunjal et al., 2024; Chen et al., 2024; Liu et al., 2024b; Chen et al., 2025a). Object hallucination refers to the phenomenon where the model can generate continuous responses, but these responses do not match the real objects within the given image. Object hallucination can undermine the capabilities of LVLMs and reduce their credibility. The researchers discovered that object hallucinations in LVLMs predominantly

arise from linguistic priors and data biases.

In this study, we conducted an in-depth analysis of the impact of LVLMs' neglect of high-frequency and low-frequency information in visual inputs on the two main causes of object hallucination in LVLMs, namely statistical bias and linguistic priors. Based on the above analysis, we propose the Multi-Frequency Contrastive Decoding (MFCD) method. The principle of MFCD is to contrast the original output distribution with the hallucination distributions generated by the input image with low-frequency information removed and the input image with high-frequency information removed, respectively. MFCD can significantly alleviate issues such as language priors and biases in pretrained data that arise due to LVLMs neglecting high-frequency or low-frequency information of image as shown in Figure 1. As a result, MFCD can effectively mitigate object hallucination in LVLMs.

Our experiments show that our proposed MFCD significantly improves the scores on object hallucination test sets such as POPE (Li et al., 2023c), CHAIR (Rohrbach et al., 2018) and MME (Fu et al., 2024), when applied to models like LLaVA 1.5 (Liu et al., 2023), LLaVA-NeXT (Li et al., 2024), and Qwen2.5-VL (Bai et al., 2025).

In summary, our main contributions are as follows:

- We, for the first time, explore the causes of object hallucinations from the perspective of the frequency domain and conduct an in-depth analysis of the impact of the lack of image frequency domain features on object hallucinations in LVLMs.

- Inspired by the above analysis, we propose a Multi-Frequency Contrastive Decoding (MFCD) method. The principle of MFCD is to contrast the original output distribution with the hallucination distributions generated by the input image with low-frequency information removed and the input image with high-frequency information removed, respectively. Our MFCD is a plug-and-play technique that requires no training and can effectively mitigate object hallucinations in LVLMs.

- Through comprehensive experiments, we demonstrated the effectiveness of the proposed MFCD in alleviating object hallucinations. Moreover, MFCD can be applied to various LVLMs without modifying the model structure or additional training, which proves the generality and robustness of the MFCD method.

## 2 Related Work

### 2.1 Large Visual-Language Models

LVLMs typically consists of a visual encoder (Dosovitskiy et al., 2021; Dehghani et al., 2023), an LLM such as LLaMA (Touvron et al., 2023) or Qwen (Yang et al., 2024; Qwen et al., 2025), and a multimodal alignment module composed of a fully connected network or a Q-Former network (Li et al., 2023a). By integrating user instructions and visual inputs, they can understand and generate diverse content in a more comprehensive way.

### 2.2 Visual Feature in Frequency Domain

The frequency domain usually reflects the intensity of color changes in an image. *High-frequency* information can clearly highlight the edges of objects. In an image, the pixel values at the edges of objects often change abruptly, and this abrupt change corresponds to the high-frequency components in the frequency domain. *Low-frequency* information, on the other hand, is mainly used to outline the general contours of objects. The low-frequency components mainly represent the areas where the pixel values change more smoothly, and these areas together form the basic shape and overall structure of the objects. In image representation learning (Li et al., 2015; Xu et al., 2020), the frequency domain information can be obtained by performing a fourier transform (Young and van Vliet, 1995; Charalampidis, 2016) on the image.

### 2.3 Hallucinations in LVLMs

Hallucinations in LVLMs refer to the contradictions between the visual inputs (regarded as facts) and the text outputs of the LVLMs (Liu et al., 2024b). Hallucinations in LVLMs may originate from data biases, the limited visual resolution of the visual encoders, misalignment in multimodal alignment module (Zhao et al., 2024), and the language priors from the internal LLMs.

Object hallucination means that the responses generated by the model do not conform to the objects in the picture (Biten et al., 2022). To evaluate the object hallucination in LVLMs, researchers have constructed a variety of evaluation datasets, such as POPE (Li et al., 2023c) and MME (Fu et al., 2024).

Current methods for mitigating hallucinations include finetuning for specific models (Liu et al., 2024a; Xing et al., 2024; Li et al., 2025), constructing preference datasets for reinforcement learning

(Sun et al., 2024; Ouali et al., 2025), or finetuning models to post-hoc rectify object hallucination in LVLMs (Zhou et al., 2024). However, these methods are usually inefficient, costly, and limited by training data and model biases. In contrast, we propose the MFCD method, which is a plug-and-play approach that does not require training. By contrast the output from the original image with the outputs from the image with low-frequency information removed and the image with high-frequency information removed respectively, it alleviates the language priors and statistical biases in LVLMs.

## 3 Method

### 3.1 Decoding of Large Visual-Language Models

The decoding process of LVLMs involves the model parameters $M$, the visual input $V$, the text query $X$, the time step $t$, and the text response $Y$. This process can be described as a series of selections made according to a strategy from the probability distribution of the model, thus generating a token sequence $Y$. This process can be expressed as:

$$Y_t \sim P(Y_t|V, X, Y_{<t}; M) \\ \propto \exp logit(Y_t|V, X, Y_{<t}; M)$$

(1)

where $Y_t$ represents the token at time step $t$, and $Y_{<t}$ represents all the tokens up to time step $(t-1)$ (Holtzman et al., 2020).

### 3.2 High-Pass Filter and Low-Pass Filter

Given the image $V \in \mathbb{R}^{H \times W \times C}$, where $H$ is the height of the image, $W$ is the width of the image, and $C$ is the number of channels of the image. The first step of both high-pass filtering and low-pass filtering is to transform the image into the frequency domain through Fourier transform:

$$F_i(u, v) = FFT(f_i(u, v)), \\ i \in \{0, 1, ..., C-1\}$$

(2)

where $F_i$ represents the frequency information in channel $i$ of the image $V$ and $FFT$ represents Fast Fourier Transform and $f_i$ represents the pixel value in channel $i$ of the image $V$. Then, we define the transfer function value $H_i^h$ of channel $i$ in high-pass filter and the transfer function value $H_i^l$ of channel $i$ in low-pass filter:

$$\begin{cases} H_i^h(u, v) = \exp\left(-\frac{D^2(u,v)}{2\sigma_h^2}\right) \\ H_i^l(u, v) = 1 - \exp\left(-\frac{D^2(u,v)}{2\sigma_h^2}\right) \end{cases}$$

(3)

where $D(u, v)$ represents the distance from the pixel point(u, v) to the center of image $V$, and $\sigma$ denotes the cutoff frequency.Then, apply filter to channel $i$ of image $V$:

$$\begin{cases} F_i^h(u, v) = Fi(u, v) \cdot H_i^h(u, v) \\ F_i^l(u, v) = Fi(u, v) \cdot H_i^l(u, v) \end{cases}$$

(4)

where $F_i^h$ represents the frequency information of channel $i$ after high-pass filtering and $F_i^l$ represents the frequency information of channel $i$ after low-pass filtering. Finally, convert the frequency domain information into an image:

$$\begin{cases} f_i^l(u, v) = FFT^{-1}(F_i^l(u, v)) \\ f_i^h(u, v) = FFT^{-1}(F_i^h(u, v)) \end{cases}$$

(5)

where $f_i^h$ is the pixel value of channel $i$ of the image $V_H$ that has undergone high-pass filtering and $f_i^l$ is the pixel value of channel $i$ of the image $V_L$ that has undergone low-pass filtering.

### 3.3 Neglect of Special Frequency Information Amplifies Object Hallucination

In order to assess the impact of neglecting frequency information on object hallucination in LVLMs, we conduct evaluations on LLaVA-1.5 using the POPE dataset and CHAIR metric (Detailed information about the POPE dataset and CHAIR evaluation metric can be found in Section 4.2.) based on MSCOCO (Lin et al., 2014), under the original image setting, "High - Pass" setting, and "Low - Pass" setting respectively. Specifically, under the "High-Pass" setting, we remove the 10% of frequency information with the lowest frequencies in the images, and under the "Low-Pass" setting, we remove the 50% of frequency information with the highest frequencies in the images.

| Visual Input | CHAIR$_i$↓ | CHAIR$_s$↓ |
|---|---|---|
| Original | 19.0 | 57.2 |
| High-Pass | 25.0 | 58.2 |
| Low-Pass | 19.9 | 57.4 |

Table 1: Result on CHAIR metric. Compared with the "Original" setting, both the CHAIR$_i$ and CHAIR$_s$ scores under the "High-Pass" and "Low-Pass" settings are improved. ↓ means "lower is better".

The evaluation results of the CHAIR metric are shown in Table 1, the absence of high-frequency or low-frequency information in the visual input leads to an increase in the values of both the CHAIR$_i$ and
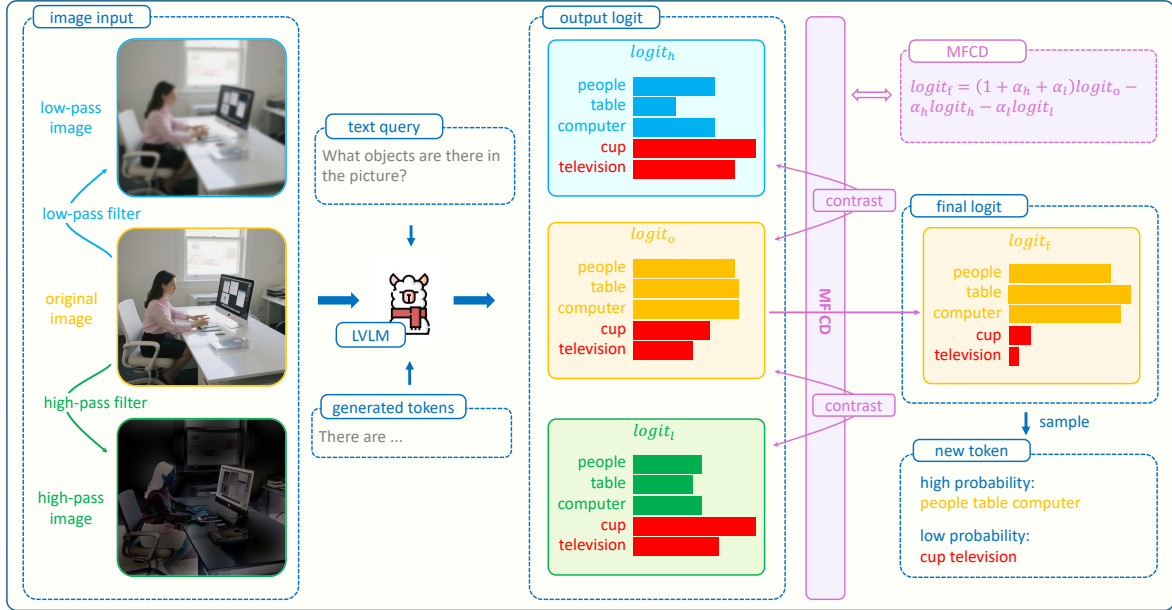
Figure 2: An illustration of Multi-Frequency Contrastive Decoding (MFCD). Hallucinatory objects are marked in red. After removing high-frequency or low-frequency image features, LVLMs are more inclined to output hallucinatory objects than when the original images are input. MFCD dynamically compares these output distributions to reduce the possibility of LVLMs outputting hallucinatory objects.

$CHAIR_s$ metrics for the model, which indicates that the likelihood of the model generating object hallucinations has increased. Additionally, we conducted a thorough case study, and the cases visually demonstrate the impact of neglecting frequency information on object hallucination in LVLMs. For more details, please refer to Appendix B.

| Visual Input | POPE | Accuracy↑ | F1 Score↑ |
|---|---|---|---|
| | Random | 81.07 | 81.88 |
| Original | Popular | 79.57 | 80.66 |
| | Adversarial | 74.50 | 76.92 |
| | Random | 74.23 | 74.26 |
| High-Pass | Popular | 72.37 | 73.25 |
| | Adversarial | 68.13 | 70.18 |
| | Random | 80.67 | 81.34 |
| Low-Pass | Popular | 77.97 | 79.44 |
| | Adversarial | 73.83 | 76.49 |

Table 2: Result on POPE dataset. Compared with the "Original" setting, both accuracy and F1 score under the "High-Pass" and "Low-Pass" settings are reduced. ↑ means "higher is better".

The evaluation results of the POPE dataset are shown in Table 2, the lack of either high-frequency or low-frequency information of visual input gives rise to a decline in both the evaluation accuracy and F1 scores of the model across the three modes of Random, Popular, and Adversarial within the POPE dataset. This indicates an elevated likelihood

of the model producing object hallucinations.

## 3.4 Multi-Frequency Contrastive Decoding

Our analysis shows that the absence of high-frequency information or low-frequency information in the image will exacerbate the language priors and statistical biases of LVLMs, thus intensifying the problem of object hallucinations in LVLMs. Inspired by the mitigation of the hallucinations of LVLMs through VCD (Leng et al., 2024), we hypothesize that separating the probability of object hallucinations, which is caused by LVLMs neglecting the high-frequency or low-frequency information in the image, from the original probability distribution may reduce object hallucinations. Based on this insight, we introduce MFCD. As shown in Figure 2, MFCD is designed to counteract the language priors and statistical biases in LVLMs by respectively contrasting the output from the original image with the outputs from the image with low-frequency information removed and the image with high-frequency information removed.

**Contrasting the Predictions** Given the model parameters $M$, the visual input $V$, the text query $X$, the time step $t$, and the text response $Y$, the visually input $V_H$ after high-pass filtering (with the low-frequency information of the $C_H$ ratios in the image removed) and the visually input $V_L$ after low-

pass filtering (with the high-frequency information of the $C_L$ ratios in the image removed), then the decoding process of MFCD can be expressed as:

$$Y_t \sim softmax[$$
$$(1 + \alpha_H + \alpha_L)logit(Y_t|V, X, Y_{<t}; M)$$
$$- \alpha_H logit(Y_t|V_H, X, Y_{<t}; M) \quad (6)$$
$$- \alpha_L logit(Y_t|V_L, X, Y_{<t}; M)]$$

where the $Y_t$ represents the token at time step $t$, the $Y_{<t}$ represents all the tokens up to time step $t-1$. The $\alpha_H$ and the $\alpha_L$ are preset hyperparameters of MFCD. The larger $\alpha_H$ is, the more intense the contrast between the output generated by the original image $V$ and the output generated by the image input $V_H$ after the low-frequency information is removed. Similarly, the larger $\alpha_L$ is, the more intense the contrast between the output generated by the original image $V$ and the output generated by the image input $V_L$ after the high-frequency information is removed.

**Adaptive Plausibility Constrains** In the MFCD method described in Formula 6, there is a possibility of wrongly penalizing reasonable tokens or wrongly rewarding unreasonable ones. To address this issue, we draw inspiration from the adaptive plausibility constraints method used in open-ended text generation (Li et al., 2023b) and add adaptive plausibility constraints to the MFCD method:

$$\mathcal{V}_{head}^{(t)} = \{Y_t \in \mathcal{V} :$$
$$P(Y_t|V_L, X, Y_{<t}; M) \quad (7)$$
$$\geq \beta \max_{\omega} P(\omega|V_L, X, Y_{<t}; M)\}$$

where $\mathcal{V}$ is the vocabulary of LVLMs and $\beta$ is a hyperparameter in $[0, 1]$, the larger the value of $\beta$ is, the more aggressive the truncation will be. And the MFCD added adaptive plausibility constraints can be expressed as:

$$Y_t \sim softmax[$$
$$(1 + \alpha_H + \alpha_L)logit(Y_t|V, X, Y_{<t}; M)$$
$$- \alpha_H logit(Y_t|V_H, X, Y_{<t}; M)$$
$$- \alpha_L logit(Y_t|V_L, X, Y_{<t}; M) \quad (8)$$
$$], \, if \, Y_t \in \mathcal{V}_{head}^{(t)}$$
$$P(Y_t|V, X, Y_{<t}; M) = 0, \, otherwise$$

which eliminates the possible negative impacts of MFCD, prevents it from generating unreasonable tokens.

**Adaptive Parameters** During the generation process, the degree of hallucination within LVLMs varies at each time step $t$. Therefore, it may be inappropriate to use fixed hyperparameters at each time step $t$ in the MFCD method. As a result, we further propose an adaptive hyperparameter improvement scheme. In this scheme, we introduce MFCD-Plus, which dynamically adjusts the parameters at each time step $t$ of the MFCD method according to the similarity between the original distribution and the hallucination distribution as well as the confidence level of the original distribution. For more details, please refer to Appendix C.

## 4 Experiment

### 4.1 Experimental Settings

**Models and Baselines** In order to validate the effectiveness of our MFCD, we carried out experiments on three representative LVLMs: LLaVA-1.5 (Liu et al., 2023), LLaVA-NeXT (Li et al., 2024) and Qwen2.5-VL (Bai et al., 2025).

We made a comparison among five decoding methods, including the sampling method (Holtzman et al., 2020), Dola (Chuang et al., 2024), and the decoding strategies specific to Large Vision-Language Models such as VCD (Leng et al., 2024) and SID (Huo et al., 2025), and our MFCD. DoLa decoding is a novel decoding strategy that dynamically selects appropriate layers by contrasting the logits of different transformer layers in large language models, aiming to enhance factuality during the decoding process, thereby reducing model hallucinations, and it demonstrates effectiveness across various tasks and models. VCD and SID respectively compare the outputs generated from the original input with the outputs generated from the input with gaussian noise added to the image and the input with visual tokens of high attention scores removed, in order to reduce hallucinations in LVLMs.

**Evaluation Settings** In sample decoding, we set the temperature to 1.2, the top-p value to 1.0, and the top-k value to 50. When evaluating the POPE dataset, both $\alpha_H$ and $\alpha_L$ of the MFCD method are set to 1.0. When evaluating the MME dataset, $\alpha_H$ of the MFCD method is set to 1.0, $\alpha_L$ is set to 0.5, and $\beta$ of the MFCD method is always set to 0.3. Moreover, in all evaluations of the MFCD method, both $C_H$ and $C_L$ are set to 0.1. In DoLa decoding, we set the candidate layers as "low", which means the first half of the Trans-

| Setting | | Random | | Popular | | Adversarial | |
|---|---|---|---|---|---|---|---|
| **Model** | **Decoding** | **Accuracy↑** | **F1 Score↑** | **Accuracy↑** | **F1 Score↑** | **Accuracy↑** | **F1 Score↑** |
| LLaVA-1.5 | Sample (base) | 81.07 | 81.88 | 79.57 | 80.86 | 74.50 | 76.92 |
| | Dola | 83.77 | 84.47 | 80.17 | 81.53 | 75.47 | 78.29 |
| | VCD | 85.23 | 86.36 | 79.57 | 82.12 | 73.10 | 77.71 |
| | SID | 86.83 | 87.27 | 82.73 | 84.04 | 75.73 | 78.74 |
| | **MFCD** | **87.07** | **87.73** | **83.07** | **84.17** | **77.03** | **79.38** |
| Qwen2.5-VL | Sample (base) | 83.27 | 80.03 | 82.24 | 78.93 | 83.27 | 80.03 |
| | Dola | 81.53 | 77.46 | 81.37 | 77.43 | 81.33 | 77.47 |
| | VCD | 84.23 | 81.39 | 83.03 | 79.91 | 82.80 | 79.98 |
| | SID | 83.26 | 79.95 | 82.73 | 79.38 | 82.23 | 78.82 |
| | **MFCD** | **85.33** | **82.91** | **84.47** | **82.01** | **84.13** | **81.75** |
| LLaVA-NeXT | Sample (base) | 71.43 | 76.63 | 64.53 | 72.51 | 65.07 | 72.72 |
| | Dola | 90.60 | 90.32 | 87.30 | 87.30 | 82.67 | 83.46 |
| | VCD | 87.73 | 87.24 | 85.83 | 85.46 | 81.60 | 81.77 |
| | SID | **91.40** | 90.46 | 86.70 | 86.87 | 83.13 | 84.20 |
| | **MFCD** | **91.40** | **91.20** | **88.17** | **88.27** | **84.50** | **85.21** |

Table 3: Results on POPE. ↑ means "higher is better". Sample refers to randomly sampling from the original output distribution. Dola, VCD, SID, and MFCD refer to randomly sampling with the same settings as the "Sample" method in the output distribution after being processed by these methods. The best performances within each setting are **bolded**.

former layers are used as candidate layers. The penalty parameter is set to 0.1, and the repetition penalty is set to 1.2. In VCD and SID, parameters are the same as those in the original paper. For more details, please refer to Appendix A.

## 4.2 Evaluation Metrics

**POPE**   POPE (Li et al., 2023c) formulates the task of evaluating object hallucinations as a binary-class visual question answering task that only requires a "yes/no" response, with the question format being "Is there a <object> in the image?".There is a 50% probability that the "<object>" truly exists. In the Random setting, for the remaining 50% of the "<objects>", they are randomly selected objects that do not exist in the image. Under the Popular setting, for the remaining 50% of the "<objects>", they are objects that do not exist in the image but frequently appear in the pre-training dataset. Under the Adversarial setting, for the remaining 50% of the "<objects>", they are objects that do not exist in the image but often co-occur with the objects in the image within the pretraining dataset. The object hallucination of LVLMs is evaluated by judging the accuracy and F1-score of the results output by LVLMs. For our analysis, we use the POPE constructed from the MSCOCO (Lin et al., 2014).

**CHAIR**   CHAIR (Rohrbach et al., 2018) is specifically designed for evaluating object hallucination in the image caption task. CHAIR consists of two evaluation metrics, namely $CHAIR_i$ ($CI$)

and $CHAIR_s$ ($CS$). The formulas for these two metrics are as follows:

$$CI = \frac{|\{\text{hallucinated objects}\}|}{|\{\text{all objects mentioned}\}|}$$

$$CS = \frac{|\{\text{sentences with hallucinated object}\}|}{|\{\text{all sentences}\}|}$$

(9)

$CHAIR_i$ and $CHAIR_s$ respectively represent the proportion of the number of hallucinated objects to the number of generated objects in image caption, and the proportion of the number of sentences containing hallucinated objects to the number of generated sentences. They evaluate the object hallucination in LVLMs at the object level and the sentence level respectively. The smaller the values of $CHAIR_i$ and $CHAIR_s$ are, the less severe the object hallucination problem in LVLMs is.

In our experiment, we randomly select 500 images from MSCOCO (Lin et al., 2014) and query LLaVA-1.5 under different decoding methods with the prompt "Please describe this picture in detail.".

**MME**   MME (Fu et al., 2024) also formulates the task of evaluating object hallucinations as a binary-class visual question answering task that only requires a "yes/no" response, with the instruction consists of a concise question and "Please answer yes or no.". MME includes perception tasks and cognitive tasks. Perception tasks can be subdivided into tasks such as coarse-grained recognition, fine-grained recognition, and OCR. Cognitive tasks can be subdivided into tasks such as com-

| Decoding | Object Level | | | | | | Attribute Level | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Existence | | | Count | | | Position | | | Color | | |
| | Acc↑ | Acc+↑ | Score↑ | Acc↑ | Acc+↑ | Score↑ | Acc↑ | Acc+↑ | Score↑ | Acc↑ | Acc+↑ | Score↑ |
| Sample (base) | 71.7 | 43.3 | 115.0 | 50.0 | 26.7 | 76.7 | 66.7 | 40.0 | 106.7 | 50.0 | 20.0 | 70.0 |
| Dola | 90.0 | 80.0 | 170.0 | 50.0 | 20.0 | 70.0 | 65.0 | 40.0 | 105.0 | 51.7 | 20.0 | 71.7 |
| VCD | 81.7 | 63.3 | 145.0 | 46.7 | 16.7 | 73.4 | 65.0 | 43.3 | 108.3 | 45.0 | 20.0 | 65.0 |
| SID | 88.3 | 76.7 | 165.0 | 46.7 | 20.0 | 66.7 | 66.7 | 40.0 | 106.7 | 40.0 | 6.7 | 46.7 |
| **MFCD** | **95.0** | **90.0** | **185.0** | **60.0** | **33.3** | **93.3** | **68.3** | **50.0** | **118.3** | **53.3** | **26.7** | **80.0** |

Table 4: Results on MME. ↑ means "higher is better". Sample refers to randomly sampling from the original output distribution. Dola, VCD, SID, and MFCD refer to randomly sampling with the same settings as the "Sample" method in the output distribution after being processed by these methods. The best performances within each setting are **bolded**.

monsense reasoning, numerical calculation, text translation, and code reasoning. MME uses accuracy and accuracy+ to measure the performance of models. Among them, accuracy+ calculates the proportion of the number of images for which the model answers both questions about the same image correctly to the total number of images, which can more strictly reflect the model's comprehensive understanding of the images. The sum of the accuracy and accuracy+ scores of each sub-task is used to obtain the total scores of the perception and cognition tasks.

**Other** More experiments for evaluating the capabilities other than alleviating hallucination, please refer to Appendix F.

### 4.3 Experimental Results

**Results on POPE** The experimental results of the POPE are summarized in Table 3. These results demonstrate the effectiveness of our MFCD method on the POPE. Our MFCD method consistently outperforms the baseline methods across all LVLMs. Additionally, our method still has significant advantages over other concurrent methods.

In addition, during the process of transitioning from the random setting to the popular setting and then to the adversarial setting, there is an overall downward trend in performance. This trend validates the view that the inherent statistical biases in large vision-language models largely contribute to the problem of object hallucination. Furthermore, our MFCD method is more effective on LLaVA-next than on Qwen2.5-VL. We speculate that this might be because the visual encoder used in Qwen-2.5-VL is more capable of capturing visual information in different frequency domains compared to the visual encoder adopted by LLaVA-NeXT. For a detailed explanation of this phenomenon, please refer to the Appendix E.

**Results on MME** The hallucination evaluation subset of MME is further divided into subsets for evaluating object hallucinations at the object level and the attribute level. Among them, the subset for evaluating hallucinations at the object level includes two types of tasks: Existence and Count, and the subset for evaluating hallucinations at the attribute level includes two types of tasks: Position and Color. We used the MME hallucination subset to evaluate various decoding methods on LLaVA-1.5, and the evaluation results are shown in Table 4. Our MFCD method can improve the performance of LVLMs at all levels of hallucination evaluation, and the improvement margin exceeds that of other decoding methods in the same period. In addition, among the four types of tasks, Count and Color perform poorly among all types of decoding methods. However, our MFCD method still has improvements in the Count and Color tasks, highlighting the advantage of the MFCD method in alleviating object hallucinations. For more detailed experimental results on MME, please refer to the Appendix D.

| Decoding | CHAIR$_i$↓ | CHAIR$_s$↓ |
|---|---|---|
| Sample (base) | 19.0 | 57.2 |
| Dola | 17.9 | 57.0 |
| VCD | 18.4 | 55.6 |
| SID | 16.2 | 54.2 |
| **MFCD** | **15.0** | **54.0** |

Table 5: Result on CHAIR. ↓ means "lower is better". The best performances within each setting are **bolded**.

**Results on CHAIR** The results on CHAIR are shown in Table 5. Compared with the baselines and all other decoding methods in the experiment, both the CHAIR$_i$ and CHAIR$_s$ scores of our MFCD method have been improved. These results demonstrate the effectiveness of our MFCD

| Hyperparameters | | | Random | | Popular | | Adversarial | |
|---|---|---|---|---|---|---|---|---|
| $\alpha_L$ | $\alpha_H$ | $\beta$ | Accuracy↑ | F1↑ | Accuracy↑ | F1↑ | Accuracy↑ | F1↑ |
| 1.0 | 1.0 | 0.3 | **91.40** | **91.20** | **88.17** | **88.27** | **84.50** | **85.21** |
| 0.0 | 1.0 | 0.3 | 75.40 | 79.35 | 71.73 | 76.83 | 70.87 | 76.27 |
| 1.0 | 0.0 | 0.3 | 75.03 | 79.29 | 71.40 | 76.52 | 66.67 | 73.11 |
| 0.0 | 0.0 | 0.3 | 66.67 | 74.23 | 67.67 | 76.40 | 62.00 | 71.78 |
| 1.0 | 1.0 | 0.0 | 80.50 | 82.16 | 78.70 | 80.70 | 67.67 | 76.40 |

Table 6: The results of the ablation experiment. ↑ means "higher is better". Moreover, $\alpha_H$, $\alpha_L$ and $\beta$ mean hyperparameters in MFCD method. The best performances within each setting are **bolded**.

method in alleviating object hallucination in the image captioning task.

**Case Study on CHAIR**     The case study shown in Appendix (Figure 5) indicates that, compared with the Sample method, the MFCD method effectively reduces the number of hallucinated objects in image captions and significantly improves the quality of the generated image captions. Besides, in this case, compared with the Sample method, the SID method fails to completely remove the hallucinated objects in the image captions. The Dola and VCD methods even introduce new hallucinated objects. In contrast, our MFCD method not only completely eliminates the hallucinated objects in the image captions but also avoids introducing new ones, demonstrating the superiority of our MFCD method.

### 4.4 Ablation Analyses

In this experiment, we conduct ablation studies to explore the roles of different modules in MFCD.

As shown in Table 6, we used the POPE dataset on LLaVA-NeXT to evaluate the performance of MFCD (with $\alpha_L$ set to 1.0, $\alpha_H$ set to 1.0 and $\beta$ set to 0.3) and **four variants of MFCD**. These four variants are: (1) removed the contrast module for the absence of *high-frequency* image information (with $\alpha_L$ set to 0.0, $\alpha_H$ set to 1.0 and $\beta$ set to 0.3) based on MFCD; (2) removed the contrast module for the absence of *low-frequency* image information (with $\alpha_L$ set to 1.0, $\alpha_H$ set to 0.0 and $\beta$ set to 0.3); (3) removed *both* the contrast modules for the absence of *high-frequency and low-frequency* image information (with $\alpha_L$ set to 0.0, $\alpha_H$ set to 0.0 and $\beta$ set to 0.3); and (4) removed the *adaptive rationality constraint* module (with $\alpha_L$ set to 1.0, $\alpha_H$ set to 1.0 and $\beta$ set to 0.0).

Compared with the MFCD methods with either the high-frequency image information absence contrast module or the low-frequency image information absence contrast module removed, the original MFCD method shows improvements in both accuracy and F1 scores on the three subsets of POPE, namely Random, Popular, and Adversarial and the MFCD method with both the high-frequency and low-frequency image information absence contrast modules removed experiences a decline in both accuracy and F1 scores on all subsets of POPE. This indicates that the high-frequency image information absence contrast module and the low-frequency image information absence contrast module can effectively alleviate the object hallucination problems in LVLMs caused by the lack of high-frequency and low-frequency information respectively.

Compared with the MFCD variant without the adaptive plausibility constrains module, the original MFCD method has improved both the accuracy and F1 scores on the three subsets of POPE, namely Random, Popular, and Adversarial. This indicates that the adaptive rationality constraint module can effectively reduce the possibility of wrongly penalizing reasonable tokens or wrongly rewarding unreasonable ones.

## 5 Conclusion

In this paper, we conduct an in-depth analysis of how the neglect of high-frequency and low-frequency information in images affects object hallucinations in LVLMs. Based on this, we propose the MFCD method, which requires no training and designed to counteract the language priors and statistical biases in LVLMs which are caused by neglecting the high-frequency or low-frequency information in the image inputs by respectively contrasting the original output distribution with the hallucination distributions generated by the input image with low-frequency information removed and the input image with high-frequency information removed. Through comprehensive experiments on multiple benchmarks across different LVLMs,

the results demonstrate that our proposed MFCD method can effectively alleviate the object hallucinations in LVLMs. Besides, we improve MFCD by introducing adaptive parameters to form MFCD-Plus. And we demonstrate through experiments that MFCD-Plus achieves performance improvement compared to MFCD.

## Limitations

Firstly, although the MFCD method has demonstrated potential in reducing object hallucination, it can only alleviate the hallucination problem caused by the neglect of high-frequency and low-frequency information in images. However, object hallucination may stem from other factors, such as the neglect of image information by LVLMs (Manevich and Tsarfaty, 2024) or incorrect user instructions (Wang et al., 2024b). Secondly, the MFCD method only supports LVLMs that process image and text inputs, and does not currently support LVLMs that handle video and text information (Zhang et al., 2024). Finally, the MFCD method may amplify other unknown biases due to contrastive decoding.

## References

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, and 8 others. 2025. Qwen2.5-vl technical report. *Preprint*, arXiv:2502.13923.

Ali Furkan Biten, Lluís Gómez, and Dimosthenis Karatzas. 2022. Let there be a clock on the beach: Reducing object hallucination in image captioning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 1381–1390.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Dimitrios Charalampidis. 2016. Recursive implementation of the gaussian filter using truncated cosine functions. *IEEE Transactions on Signal Processing*, 64(14):3554–3565.

Guoqing Chen, Fu Zhang, Jinghao Lin, Chenglong Lu, and Jingwei Cheng. 2025a. Rrhf-v: Ranking responses to mitigate hallucinations in multimodal large language models with human feedback. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 6798–6815.

Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2025b. Sharegpt4v: Improving large multi-modal models with better captions. In *Computer Vision – ECCV 2024*, pages 370–387, Cham. Springer Nature Switzerland.

Xuweiyi Chen, Ziqiao Ma, Xuejun Zhang, Sihan Xu, Jianing Yang, David F. Fouhey, Joyce Chai, and Shengyi Qian. 2024. Multi-object hallucination in vision language models. In *Advances in Neural Information Processing Systems*, volume 37, pages 44393–44418. Curran Associates, Inc.

Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R. Glass, and Pengcheng He. 2024. Dola: Decoding by contrasting layers improves factuality in large language models. In *The Twelfth International Conference on Learning Representations*.

Mostafa Dehghani, Basil Mustafa, Josip Djolonga, Jonathan Heek, Matthias Minderer, Mathilde Caron, Andreas Steiner, Joan Puigcerver, Robert Geirhos, Ibrahim M Alabdulmohsin, Avital Oliver, Piotr Padlewski, Alexey Gritsenko, Mario Lucic, and Neil Houlsby. 2023. Patch n' pack: Navit, a vision transformer for any aspect ratio and resolution. In *Advances in Neural Information Processing Systems*, volume 36, pages 2252–2274. Curran Associates, Inc.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*.

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. 2024. Mme: A comprehensive evaluation benchmark for multimodal large language models. *Preprint*, arXiv:2306.13394.

Zigang Geng, Binxin Yang, Tiankai Hang, Chen Li, Shuyang Gu, Ting Zhang, Jianmin Bao, Zheng

Zhang, Houqiang Li, Han Hu, Dong Chen, and Baining Guo. 2024. Instructdiffusion: A generalist modeling interface for vision tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12709–12720.

Anisha Gunjal, Jihan Yin, and Erhan Bas. 2024. Detecting and preventing hallucinations in large vision language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):18135–18143.

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*.

Fushuo Huo, Wenchao Xu, Zhong Zhang, Haozhao Wang, Zhicheng Chen, and Peilin Zhao. 2025. Self-introspective decoding: Alleviating hallucinations for large vision-language models. In *The Thirteenth International Conference on Learning Representations*.

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 *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13872–13882, Los Alamitos, CA, USA. IEEE Computer Society.

Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. 2024. Llava-next: Stronger llms supercharge multimodal capabilities in the wild (2024). *URL https://llava-vl. github. io/blog/2024-05-10-llava-next-stronger-llms*.

Jia Li, Ling-Yu Duan, Xiaowu Chen, Tiejun Huang, and Yonghong Tian. 2015. Finding the Secret of Image Saliency in the Frequency Domain . *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 37(12):2428–2440.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. BLIP-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 19730–19742. PMLR.

Shuo Li, Jiajun Sun, Guodong Zheng, Xiaoran Fan, Yujiong Shen, Yi Lu, Zhiheng Xi, Yuming Yang, Wenming Tan, Tao Ji, Tao Gui, Qi Zhang, and Xuanjing Huang. 2025. Mitigating object hallucinations in mllms via multi-frequency perturbations. *Preprint*, arXiv:2503.14895.

Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023b. Contrastive decoding: Open-ended text generation as optimization. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12286–12312, Toronto, Canada. Association for Computational Linguistics.

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

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*, pages 740–755, Cham. Springer International Publishing.

Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2024a. 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. 2024b. A survey on hallucination in large vision-language models. *Preprint*, arXiv:2402.00253.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. In *Advances in Neural Information Processing Systems*, volume 36, pages 34892–34916. Curran Associates, Inc.

Ana P Majtey, Pedro W Lamberti, and Domingo P Prato. 2005. Jensen-shannon divergence as a measure of distinguishability between mixed quantum states. *Physical Review A—Atomic, Molecular, and Optical Physics*, 72(5):052310.

Avshalom Manevich and Reut Tsarfaty. 2024. Mitigating hallucinations in large vision-language models (LVLMs) via language-contrastive decoding (LCD). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 6008–6022, Bangkok, Thailand. Association for Computational Linguistics.

Yassine Ouali, Adrian Bulat, Brais Martinez, and Georgios Tzimiropoulos. 2025. Clip-dpo: Vision-language models as a source of preference for fixing hallucinations in lvlms. In *Computer Vision – ECCV 2024*, pages 395–413, Cham. Springer Nature Switzerland.

Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, and 25 others. 2025. Qwen2.5 technical report. *Preprint*, arXiv:2412.15115.

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. 2021. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of

*Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.

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, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *Preprint*, arXiv:2302.13971.

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *Preprint*, arXiv:2409.12191.

Xintong Wang, Jingheng Pan, Liang Ding, and Chris Biemann. 2024b. Mitigating hallucinations in large vision-language models with instruction contrastive decoding. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 15840–15853, Bangkok, Thailand. Association for Computational Linguistics.

Shangyu Xing, Fei Zhao, Zhen Wu, Tuo An, Weihao Chen, Chunhui Li, Jianbing Zhang, and Xinyu Dai. 2024. EFUF: Efficient fine-grained unlearning framework for mitigating hallucinations in multimodal large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1167–1181, Miami, Florida, USA. Association for Computational Linguistics.

Kai Xu, Minghai Qin, Fei Sun, Yuhao Wang, Yen-Kuang Chen, and Fengbo Ren. 2020. Learning in the Frequency Domain . In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1737–1746, Los Alamitos, CA, USA. IEEE Computer Society.

Shiyue Xu, Fu Zhang, Jingwei Cheng, and Linfeng Zhou. 2025. Mwpo: Enhancing llms performance through multi-weight preference strength and length optimization. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 20566–20581.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 43 others. 2024. Qwen2 technical report. *CoRR*, abs/2407.10671.

Gokul Yenduri, Ramalingam M, Chemmalar Selvi G, Supriya Y, Gautam Srivastava, Praveen Kumar Reddy Maddikunta, Deepti Raj G, Rutvij H Jhaveri, Prabadevi B, Weizheng Wang, Athanasios V. Vasilakos, and Thippa Reddy Gadekallu. 2023. Generative pre-trained transformer: A comprehensive review on enabling technologies, potential applications, emerging challenges, and future directions. *Preprint*, arXiv:2305.10435.

Ian T. Young and Lucas J. van Vliet. 1995. Recursive implementation of the gaussian filter. *Signal Processing*, 44(2):139–151.

Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. 2024. Video instruction tuning with synthetic data. *Preprint*, arXiv:2410.02713.

Zhiyuan Zhao, Bin Wang, Linke Ouyang, Xiaoyi Dong, Jiaqi Wang, and Conghui He. 2024. Beyond hallucinations: Enhancing lvlms through hallucination-aware direct preference optimization. *Preprint*, arXiv:2311.16839.

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

## A   Detailed Experimental Setup

For the POPE experiment and the MME experiment, we set the maximum number of tokens to be generated as 2048. In the image caption task for evaluating the CHAIR metric, to prevent the generated image captions from being truncated, we set the maximum number of tokens to be generated as 4096.

## B   Neglect of Special Frequency Information Amplifies Object Hallucination

As shown in the cases in Figure 3 and Figure 4, after removing the high-frequency information (under the "Low-Pass" setting) and the low-frequency information (under the "High-Pass" setting) from the image input, the number of hallucinated objects in the image captions generated by LVLMs has significantly increased. This further proves that neglect of high-frequency information or low-frequency information in image input amplifies object hallucination in image caption task. Besides, under the "High-Pass" setting where the low-frequency information has been removed, LVLMs will recognize the scene in the image as a black room. This is probably because the low-frequency information represents the image information of the areas with gentle color changes in the image, and the removal of the low-frequency information leads to the background color in the image turning black.

## C   Adaptive Parameters

We propose MFCD-Plus. In MFCD-Plus, we use the Jensen-Shannon divergence (Majtey et al., 2005) to measure the similarity between the original distribution and the hallucination distribution. The higher the similarity between the original distribution and the hallucination distribution is, the more severe the relevant hallucination in the LVLMs is, and a larger contrast parameter is needed. The formula related to adjusting the contrast parameter is as follows:

$$
\begin{aligned}
\alpha_H^{(t)} = \alpha_H / JSD( \\
softmax(logit(Y_t|V, X, Y_{<t}; M)), \\
softmax(logit(Y_t|V_H, X, Y_{<t}; M))) \\
\alpha_L^{(t)} = \alpha_L / JSD( \\
softmax(logit(Y_t|V, X, Y_{<t}; M)), \\
softmax(logit(Y_t|V_L, X, Y_{<t}; M)))
\end{aligned} \quad (10)
$$

where JSD is the calculation formula of the Jensen-Shannon divergence. And $\alpha_H^{(t)}$ and $\alpha_L^{(t)}$ respectively control the degree of contrast between the original distribution in the MFCD method and the hallucination distribution caused by ignoring low-frequency information, as well as the degree of contrast between the original distribution and the hallucination distribution caused by ignoring high-frequency information at time step $t$.

In addition, we use conditional entropy to measure the degree of uncertainty within LVLMs. The higher the degree of uncertainty, the greater the need for plausibility constrains. The formula for adjusting the parameter in plausibility constrainsis as follows:

$$
\beta^{(t)} = \beta \times (1 - e^{-H_{LVLM}(Y_t|V,X,Y_{<t};M)}) \quad (11)
$$

where $H_{LVLM}(Y_t|V, X, Y_{<t}; M)$ represents the conditional entropy of the original distribution generated by LVLMs at time step $t$, and $\beta^{(t)}$ is used to control the intensity of the plausibility constraints at time step $t$.

In conclusion, MFCD-Plus can be described by the following formula:

$$
\begin{aligned}
Y_t \sim softmax[ \\
(1 + \alpha_H^{(t)} + \alpha_L^{(t)})logit(Y_t|V, X, Y_{<t}; M) \\
- \alpha_H^{(t)}logit(Y_t|V_H, X, Y_{<t}; M) \\
- \alpha_L^{(t)}logit(Y_t|V_L, X, Y_{<t}; M) \\
], \; if \; Y_t \in \mathcal{V}_{head}^{(t)} \\
P(Y_t|V, X, Y_{<t}; M) = 0, \; otherwise
\end{aligned} \quad (12)
$$

where $\mathcal{V}_{head}^{(t)}$ can be described by the following formula:

$$
\begin{aligned}
\mathcal{V}_{head}^{(t)} = \{Y_t \in \mathcal{V} : \\
P(Y_t|V_L, X, Y_{<t}; M) \\
\geq \beta^{(t)} \max_{\omega} P(\omega|V_L, X, Y_{<t}; M)\}
\end{aligned} \quad (13)
$$

We evaluated the MFCD-Plus method on LLaVA-1.5 using the POPE dataset. As shown in Table 7, compared with the MFCD method, the MFCD-Plus method has improved accuracy and F1 scores under the three settings of Random, Popular, and Adversarial, which reflects the effectiveness of the improvement of the adaptive parameters.

## D   Detailed Experimental Results on MME

We evaluated the perception tasks of MME on the LLaVA platform, and the evaluation results (Score) are presented in the Table 8.

| Decoding | Random | | Popular | | Adversarial | |
|---|---|---|---|---|---|---|
| | **Accuracy**↑ | **F1 Score**↑ | **Accuracy**↑ | **F1 Score**↑ | **Accuracy**↑ | **F1 Score**↑ |
| **MFCD** | 87.07 | 87.73 | 83.07 | 84.17 | 77.03 | 79.38 |
| **MFCD-Plus** | **88.53** | **88.66** | **84.47** | **85.28** | **77.43** | **79.85** |

Table 7: Results on POPE. ↑ means "higher is better". The best performances within each setting are **bolded**.

| Decoding | Existence | Count | Position | Color | OCR | Posters | Celebrity | Scene | Landmark | Artwork |
|---|---|---|---|---|---|---|---|---|---|---|
| Sample (base) | 115.0 | 76.7 | 106.7 | 70.0 | **132.5** | 106.8 | 95.6 | 127.5 | 87.5 | 87.3 |
| Dola | 170.0 | 70.0 | 105.0 | 71.7 | 107.5 | 113.5 | 110.3 | 138.8 | 110.5 | 99.0 |
| VCD | 145.0 | 73.4 | 108.3 | 65.0 | 127.5 | 127.9 | 118.2 | 139.8 | 114.0 | 94.0 |
| SID | 165.0 | 66.7 | 106.7 | 46.7 | 100.0 | 124.1 | 118.5 | **160.5** | 105.8 | 104.2 |
| **MFCD** | **185.0** | **93.3** | **118.3** | **80.0** | 117.5 | **133.3** | **143.5** | 158.5 | **133.5** | **111.0** |

Table 8: Results on perception tasks of MME. Sample refers to randomly sampling from the original output distribution. Dola, VCD, SID, and MFCD refer to randomly sampling with the same settings as the "Sample" method in the output distribution after being processed by these methods. The best performances within each setting are **bolded**.

Our method, MFCD, achieves the best performance on 8 out of the 10 subsets, which further demonstrates the effectiveness of our approach.

# E Explanation of Experimental Results on POPE

We further analyzed the sensitivity of LLaVA-NeXT and Qwen2.5-VL to frequency-domain information in images. We start with low frequencies to analyze the sensitivity of Qwen2.5-VL and LLaVA-NEXT to frequency-domain information. To this end, we designed the following experiments: we evaluated the POPE metric on LLaVA-NeXT and Qwen2.5-VL using both original images and images with 10% low-frequency information removed (high-Pass images), respectively.

The experimental results are shown in Table 9. For Qwen2.5-VL, when using images with 10% low-frequency information removed as inputs compared to original images, both the accuracy and F1 score in the POPE metric showed a significant decline. In contrast, for LLaVA-NeXT, the decreases in accuracy and F1 score of the POPE metric were relatively small under the same condition.

From the experimental results, we can infer that low-frequency information in images is critical for Qwen2.5-VL, as the model's generation process highly relies on such information. Therefore, removing 10% of the low-frequency information leads to a substantial drop in the POPE metric for Qwen2.5-VL. By contrast, LLaVA-NeXT does not make sufficient use of low-frequency information in images, and its generation process hardly depends on such information. As a result, after re-

moving low-frequency information from images, the POPE metric of LLaVA-NeXT only decreases slightly.

From the above analysis, it can be seen that Qwen2.5-VL makes more sufficient use of low-frequency information. Therefore, it can be inferred that object hallucinations caused by the lack of frequency-domain information rarely occur in Qwen2.5-VL, so MFCD has a limited helpful effect on Qwen2.5-VL. Similarly, LLaVA-NeXT has a low utilization rate of low-frequency information, so LLaVA-NeXT is more prone to object hallucinations caused by the lack of frequency-domain information. Therefore, the MFCD method has a significant improvement effect on object hallucinations caused by the lack of frequency-domain information.

# F Experiment for evaluating the capabilities other than alleviating hallucination

We evaluated the capabilities of the model other than alleviating hallucination on LLaVA-Bench (in-the-wild) dataset. Different from datasets like POPE and MME where models only need to answer Yes or No, LLaVA-Bench (in the wild) requires models to generate detailed responses to questions and evaluates the models' generation-related performance by scoring the answers according to certain principles.

In the process of generating the answer, in sample, we set the temperature to 1.2, the top-p value to 1.0, and the top-k value to 50, and in MFCD, both $\alpha_H$ and $\alpha_L$ of the MFCD method are set to 1.0.

| Setting | | Random | | Popular | | Adversarial | |
|---|---|---|---|---|---|---|---|
| Model | Image Input | Accuracy↑ | F1 Score↑ | Accuracy↑ | F1 Score↑ | Accuracy↑ | F1 Score↑ |
| Qwen2.5-VL | Sample (base) | 83.27 | 80.03 | 82.24 | 78.93 | 83.27 | 80.03 |
| | High-Pass | 68.33 | 55.02 | 66.43 | 53.70 | 64.63 | 51.66 |
| LLaVA-NeXT | Sample (base) | 71.43 | 76.63 | 64.53 | 72.51 | 65.07 | 72.72 |
| | High-Pass | 65.03 | 63.15 | 63.43 | 64.43 | 61.80 | 61.78 |

Table 9: Results on POPE. ↑ means "higher is better".

Subsequently, we employed DeepSeek-V3 to evaluate the generated answers and score answers across dimensions including helpfulness, relevance, accuracy, and level of detail. Next, we calculated the average scores on tasks **Conversation**, **Detail description**, and **Complex reasoning**, as well as the average score across the entire dataset, to measure the comprehensive performance of model outputs in four aspects: helpfulness, relevance, accuracy, and level of details.

| Decoding | Conversation | Detail description | Complex reasoning |
|---|---|---|---|
| Sample (base) | 4.24 | 3.27 | 4.57 |
| **MFCD** | **4.47** | **3.73** | **5.11** |

Table 10: Results on LLaVA-Bench (in the wild). Sample refers to randomly sampling from the original output distribution. MFCD refers to randomly sampling with the same settings as the "Sample" method in the output distribution after being processed by this methods. The best performances within each setting are **bolded**.

The experimental results are shown in Table 10. Our MFCD method can not only enhance the model's ability to alleviate hallucinations but also improve other capabilities of the model beyond hallucination alleviation, such as the quality of generated content.

**Prompt**

Please describe this picture in detail.

**Sampel:** The image showcases a vibrant outdoor booth selling a variety of food and drinks. Several hot dogs are hanging over a counter or display, along with multiple bananas distributed across the scene. These hot dogs and bananas are likely an important part of the stand's sales in a busy location. In the vicinity of the booth, a few people are present, perhaps enjoying the food or the general atmosphere. Furthermore, there are other individuals visible in the background, likely engaged in their daily activities or simply passing by the area.

**High-Pass:** The scene presents a dark kitchen or dining room area with multiple colored objects. There are several bowls in various shapes and sizes placed around the space. One bowl is larger and placed more prominently toward the left side; several other smaller bowls are located towards the right side of the room. In the middle of the room, there is a small bar that accommodates several knives. Additionally, there are other knives found in different positions within the room, including in the middle of the dining table. A sink may also be present within the room, but due to the dark colors and low presence of lighting, the sink cannot be confidently determined. Overall, the setting appears to be a mixture of a kitchen and dining room, with bowls dispersed throughout each area.

**Low-Pass:** The image showcases a food stand located outdoors. The food stand offers various food items, with an oven loaded up with delicious donuts visible. The delicious items, specifically the donuts and other snacks, are visible from every corner. Two different sections of hamburgers caught the attention at this food stand, with hamburgers displayed in all shapes and sizes, inviting customers to choose a flavor. In addition to the appetizing treats, several bottles and bananas can be found among the items for sale. This mouthwatering offer might be a tempting scene for visitors, especially hungry or thirsty people who are visiting the nearby event.

Figure 3: Case Study used to give an illustration of neglect of high-frequency information or low-frequency information amplifies object hallucination in image caption task. The figure contains the image caption generated by the Sample method and the image captions generated under High-Pass and Low-Pass settings. Hallucination objects are marked in red.

**Prompt**

Please describe this picture in detail.

---

**Sampel:** In the image, there are two sandwiches prepared with different types of meats. One sandwich is stuffed with ham and the other includes bacon. These sandwiches appear to be served in a restaurant or deli setting, attractively displayed in a display case. A sign reads Special La Costa Pancino for sale in the restaurant, indicating the sandwich is for sale for $6.90. Additionally, there's a bowl in the background near the left side of the screen, likely filled with other food options or condiments.

---

**High-Pass:** In the dark room, there are various items scattered on the dining room table. One main focus is a table full of food, including sandwiches and a cake. In front of the table, there are two sandwiches placed on a plate, one with a bowl near it, possibly containing snacks to accompany the sandwiches. There are also two plastic glasses and a handbag placed nearby, adding to the overall setting. These items contribute to a well-prepared dining experience despite the dim lighting.

---

**Low-Pass:** The image displays a hot lunch pantry featuring a delicious selection of sandwiches and desserts on a serving trolley. There are three sandwiches, including what appears to be Italian sandwiches. Two desserts can be seen on the server, which are cake and pizza displayed at different spots within the trolley. Additionally, there is a slice of cake placed further down below. Around the pantry is a dining table, with a second dining table close to the right corner of the space. Each dining table is adorned with several chairs surrounding them. This inviting setup is perfect for people who are hungry for such offerings as the three types of sandwiches and two cakes.

Figure 4: Case Study used to give an illustration of neglect of high-frequency information or low-frequency information amplifies object hallucination in image caption task. The figure contains the image caption generated by the Sample method and the image captions generated under High-Pass and Low-Pass settings. Hallucination objects are marked in red.

**Prompt**

Please describe this picture in detail.

**Sampel:** The image showcases a vibrant outdoor booth selling a variety of food and drinks. Several hot dogs are hanging over a counter or display, along with multiple bananas distributed across the scene. These hot dogs and bananas are likely an important part of the stand's sales in a busy location. In the vicinity of the booth, a few people are present, perhaps enjoying the food or the general atmosphere. Furthermore, there are other individuals visible in the background, likely engaged in their daily activities or simply passing by the area.

**Dola:** The image showcases a marketplace stand filled with an array of skewers and kebab food, including bananas, doughnuts, and hot dogs. The doughnuts are located towards the left side, while bananas are situated near the left edge of the image. The hot dogs are found scattered throughout the scene, with some of them placed higher than the others. In addition to the food items, there is a person standing towards the left side of the image, possibly attending to the stand or browsing the products. A backpack is visible near the top right corner of the stand, suggesting a customer may have placed their belongings nearby.

**VCD:** The image features an outdoor street vendor selling a variety of sweets, with a particular focus on bananas and sticky buns. There is both a table covered by a tablecloth filled with the desserts and a tent that extends above the table. The table is filled with bananas, possibly both yellow and ripe varieties. There are two main dining tables shown in the scene, covered with banana goods. Additionally, there are a couple of chairs visible in the background, suggesting the area where customers can sit, relax, and enjoy their treats. An umbrella can also be spotted extending above the dining area, providing shade and protection for the vendors and customers.

**SID:** The image showcases a food stand at an outdoor market, featuring a variety of food items. There are multiple bananas displayed, some of which are placed closer to the foreground, while others are scattered around the stand. A large number of hot dogs are also visible, with some placed near the bananas and others situated further away. A person is present in the scene, possibly a customer or a vendor, located near the top left corner of the image. The market stand is well-stocked with a diverse assortment of food items, creating an inviting atmosphere for visitors.

**MFCD:** The image features a street food stand selling various types of bananas. There are several bananas displayed at the stand, with some located near the left side of the scene and others towards the right side. In the background, there are a few people, with one person standing on the left side, another person on the right side, and a third person further back on the right side. Additionally, there is a backpack placed near the center of the scene.

Figure 5: Case Study from CHAIR on LLaVA-1.5. The figure contains the image captions generated by the Sample method, the Dola method, the VCD method, the SID method and the MFCD method. Hallucination objects are marked in red.