

Mitigating Hallucination in Multimodal Large Language Model via Hallucination-targeted Direct Preference Optimization

Anonymous ACL submission

Abstract

Multimodal Large Language Models (MLLMs) are known to hallucinate, which limits their practical applications. Recent works have attempted to apply Direct Preference Optimization (DPO) to enhance the performance of MLLMs, but have shown inconsistent improvements in mitigating hallucinations. To address this issue more effectively, we introduce Hallucination-targeted Direct Preference Optimization (HDPO) to reduce hallucinations in MLLMs. Unlike previous approaches, our method tackles hallucinations from their diverse forms and causes. Specifically, we develop three types of preference pair data targeting the following causes of MLLM hallucinations: (1) insufficient visual capabilities, (2) long context generation, and (3) multimodal conflicts. Experimental results demonstrate that our method achieves superior performance across multiple hallucination evaluation datasets, surpassing most state-of-the-art (SOTA) methods and highlighting the potential of our approach. Ablation studies and in-depth analyses further confirm the effectiveness of our method and suggest the potential for further improvements through scaling up.

1 Introduction

Large Language Models (LLMs) have been verified in various fields, demonstrating their potential (OpenAI, 2024; Dubey et al., 2024; Sun et al., 2024), while they encounter challenges such as hallucination. Multimodal Large Language Models (MLLMs) are also known to hallucinate (Bai et al., 2024). Specifically, they often produce unfaithful content that does not align with the visual input, which undermines their reliability and practicality, particularly in critical applications such as autonomous driving (Cui et al., 2024) or medical tasks (Liu et al., 2023a). Hence, addressing MLLM hallucination (**M-hallu**) is essential.

Recently, some pioneer preference optimization methods like Direct Preference Optimization (DPO) (Rafailov et al., 2024) have been introduced, which encourages the model to learn from comparisons between positive and negative samples, alleviating hallucinations (Zhao et al., 2023; Pi et al., 2025). However, most current methods cannot deliver consistent improvements across all types of M-hallu tasks (e.g., VQA and captioning tasks, as shown in our experiments of Table 1). Additionally, it appears that the model’s improvement on specific tasks is closely related to the format of the training data. For instance, the data of SeVa (Zhu et al., 2024) primarily consists of VQA, which explains its strong performance on VQA-related hallucination evaluation. However, its results on captioning tasks are relatively unsatisfactory. Moreover, these methods do not explicitly consider diverse sources of M-hallu. Hence, we argue that if we focus on mitigating multimodal hallucinations, we should be able to address diverse types of hallucination causes and tasks, and design hallucination-targeted preference pairs for DPO accordingly. Our goal is to comprehensively alleviate all multimodal hallucination problems, including both discriminative tasks (e.g., VQA) and generative tasks (e.g., image captioning).

Different from the hallucinations in LLMs, M-hallu primarily arises from the following three aspects: (1) **Insufficient visual capability**: This occurs when the MLLM’s visual encoder lacks the necessary strength, being distracted by relatively unimportant visual information, leading to hallucinations; (2) **Incapable long-context generation**: We observe that hallucinations become more pronounced as the generated content grows longer, similar to long-range forgetting, which needs to be addressed in practical applications; (3) **Multimodal conflicts**: Multimodal conflicts frequently arise in real-world scenarios due to the inevitable noises in texts and images. MLLMs are more prone to

082 hallucinations with conflicting information existing
083 between text and image (Liu et al., 2024c).

084 To address the aforementioned challenges, we
085 propose **Hallucination-targeted Direct Prefer-**
086 **ence Optimization (HDPO)** to mitigate M-hallu.
087 Our approach constructs hallucination-targeted
088 preference pairs, specifically designed to address
089 various forms and causes of hallucinations. Specif-
090 ically, we design three types of DPO data reflect-
091 ing the corresponding hallucination causes as fol-
092 lows: (1) For *insufficient visual capability*, during
093 the model’s autoregressive decoding, we preserve
094 only some visual tokens with the lowest attention
095 scores to produce targeted negative responses that
096 reflect incorrect visual information distraction, urg-
097 ing MLLMs to pay attention to more effective vi-
098 sual information. (2) For *incapable long context*
099 *generation*, we specifically select positive exam-
100 ples from high-quality long-form captions, while
101 creating negative examples where the latter part
102 of the response deviates from the image content,
103 simulating long-form hallucinations. (3) For *mul-*
104 *timodal conflicts*, we add conflicting information
105 with images into prompts to generate negative ex-
106 amples. We provide positive and negative pairs
107 with questions featuring conflicting prefixes to train
108 the model to correctly respond to the question even
109 containing conflicting information.

110 We conduct extensive experiments to evaluate
111 our approach across various types of M-hallu tasks.
112 The results demonstrate that our HDPO framework
113 achieves the overall best performance in effectively
114 mitigating MLLM hallucinations on various tasks.
115 Our contributions are summarized as follows:

- 116 • We analyze three key causes behind MLLM
117 hallucinations from visual capability, long-
118 context generation, and multimodal conflicts
119 aspects, offering valuable insights to guide
120 future advancements.
- 121 • Based on these analyses, we propose a novel
122 HDPO, aiming to jointly address all types of
123 M-hallu tasks. To the best of our knowledge,
124 we are the first to adopt hallucination-targeted
125 DPO from diverse aspects with our novel DPO
126 data construction strategies.
- 127 • Through extensive experiments on different
128 datasets, HDPO demonstrates consistent im-
129 provements in all types of M-hallu tasks.

2 Related Work 130

Hallucinations in MLLMs. Recently, the rapid
131 progress of LLMs has accelerated the MLLMs,
132 demonstrating impressive visual understanding
133 ability. However, they still encounter hallucina-
134 tions. Lots of works have explored various ap-
135 proaches to mitigate M-hallu. Some training-free
136 methods are proposed including enhancing mod-
137 els’ decoding process (Leng et al., 2024; Huang
138 et al., 2024; Chen et al., 2024) and utilizing exter-
139 nal feedbacks to reduce hallucinations (Yin et al.,
140 2023; Wu et al., 2024), while other training meth-
141 ods enhance datasets’ quality (Liu et al., 2023b).
142 Our work belongs training category. And we will
143 elaborate more on related preference optimizaiton
144 methods for improving MLLMs below. 145

Preference Optimization on MLLMs. Recently,
146 preference optimization like DPO has been used
147 to enhance models. HA-DPO (Zhao et al., 2023)
148 views hallucinations as models’ preferences. By
149 leveraging ChatGPT (Achiam et al., 2023) along-
150 side ground truth annotations from existing im-
151 age datasets, it generates positive examples aligned
152 with image content, while the model’s original out-
153 puts serve as negative examples for direct prefer-
154 ences optimization. Although effective, the con-
155 struction of negative examples is suboptimal, as
156 it may not fully capture the diverse forms of M-
157 hallu. SeVa (Zhu et al., 2024) generates negative
158 examples by introducing noise into images and
159 treats the model’s original outputs as positive ex-
160 amples, constructing pairs for DPO. In addition to
161 adding noise, BPO (Pi et al., 2025) injects errors
162 into positive examples via the LLM backbone of
163 MLLMs to construct negative examples. However,
164 our experiments indicate that while these methods
165 demonstrate strong capabilities, their performance
166 in hallucination-related evaluations is not particu-
167 larly impressive. Nonetheless, these works demon-
168 strate the superiority of DPO in enhancing models’
169 capabilities. Inspired by these findings, we aim to
170 develop methods to further mitigate M-hallu from
171 its diverse forms with hallucination-targeted direct
172 preference optimization. 173

HDPO differs from existing methods. Unlike
174 other existing preference optimization approaches,
175 we primarily focus on hallucination-targeted prefer-
176 ence optimization. We analyze and address halluci-
177 nations in MLLMs from diverse causes and forms.
178 During the preference optimization process, the
179 model learns to distinguish between positive and
180

negative examples. HA-DPO enables the model to be aware of hallucinated content in its original outputs, though its effectiveness is limited to capturing the diverse range of hallucinations as the data is insufficient. In contrast, other works use general preference data, which improves overall model capability but shows inconsistency across different hallucination benchmarks. Therefore, we aim to enhance the effectiveness of DPO by constructing examples that reflect a wider range of hallucination forms and characteristics, allowing the model to align better to make less hallucination.

Causes of hallucinations in MLLMs. There are substantial works exploring M-hallu, offering insightful perspectives. VCD suggests that language prior within MLLM is a key factor in inducing hallucinations. The Less is More (Yue et al., 2024) highlights that hallucinations are more prevalent in longer texts. In contrast, Eyes Wide Shut (Tong et al., 2024) identifies limitations in the current CLIP-based visual encoders used in MLLMs, showing that they fail to capture fine-grained details. Furthermore, SID (Huo et al., 2024) points out that tokens with lower weights in the early layers can trigger subsequent hallucinations. Meanwhile, PhD (Liu et al., 2024c) demonstrates that M-hallu stems from conflicts between multimodal information, and counterintuitive images particularly prone to causing hallucinations. Collectively, these studies provide valuable insights into understanding and addressing M-hallu.

3 Method

In this section, we provide a brief preliminaries of MLLM and DPO, followed by a detailed explanation of our proposed HDPO for constructing three types of hallucination-targeted preference data.

3.1 Preliminaries

Multimodal Large Language Models. MLLMs utilize LLMs to predict the probability distribution of the next token for each textual input. Given a prompt x that includes both an image and a text query, MLLMs generate a corresponding text response y . By incorporating visual information, MLLMs enhance the capabilities of LLMs, enabling multimodal understanding.

Direct Preference Optimization. To better align LLMs with human preferences, preference optimization methods have been developed. Among these, Reinforcement Learning from Human Feed-

back (RLHF) is a widely recognized method, though it involves training a reward model, which can be quite challenging. In contrast, Direct Preference Optimization (DPO) (Rafailov et al., 2024) utilizes preferences data directly, without the need for a reward model. This makes DPO the approach we employ. Given a pre-processed preference dataset D containing x , y_c , and y_r , where x represents the input prompt, y_c is the preferred response, and y_r is the rejected response, DPO optimizes the language model through the following loss function:

$$\mathcal{L}_d = -\mathbb{E}_D \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_c|x)}{\pi_{\text{ref}}(y_c|x)} - \beta \log \frac{\pi_\theta(y_r|x)}{\pi_{\text{ref}}(y_r|x)} \right) \right],$$

where $\pi_{\text{ref}}(y|x)$ denotes the reference policy, i.e., the language model after supervised fine-tuning, with θ as the trainable parameter.

Motivation of our HDPO. We propose HDPO, which constructs high-quality preference pairs related to the major causes of MLLM hallucinations with DPO to alleviate M-hallu. Note that the main contributions of HDPO lie in the discovery, analysis, and appropriate sample constructions of three representative types of M-hallu. Enhanced DPO algorithm is promising but not our focus.

3.2 Overview of HDPO

The primary goal of HDPO is to broadly tackle various M-hallu issues by constructing hallucination-targeted preference pairs, rather than relying on DPO data of specific tasks. Without loss of generality, we adopt a generalized data format: image-descriptive text data, which we believe more effectively captures various forms of hallucination.

For DPO in MLLMs, we require a preference dataset D , denoted as (I, q, y_c, y_r) , where I is the image, q is the question, y_c is the preferred (positive) response, and y_r is the rejected (negative) response. Currently, there are already many high-quality positive examples available, such as the refined positive examples in HA-DPO for the VG dataset, which leverage ChatGPT to enhance image annotations, and a vast number of positive examples labeled by GPT-4V in ShareGPT4V (Chen et al., 2023). These high-quality datasets have a strong alignment with the image content, making them suitable for use as positive examples in DPO. Therefore, our focus going forward is on how to construct more valuable and informative negative examples, particularly those that target hallucination, which will help the model learn preferences and reduce hallucination occurrences.

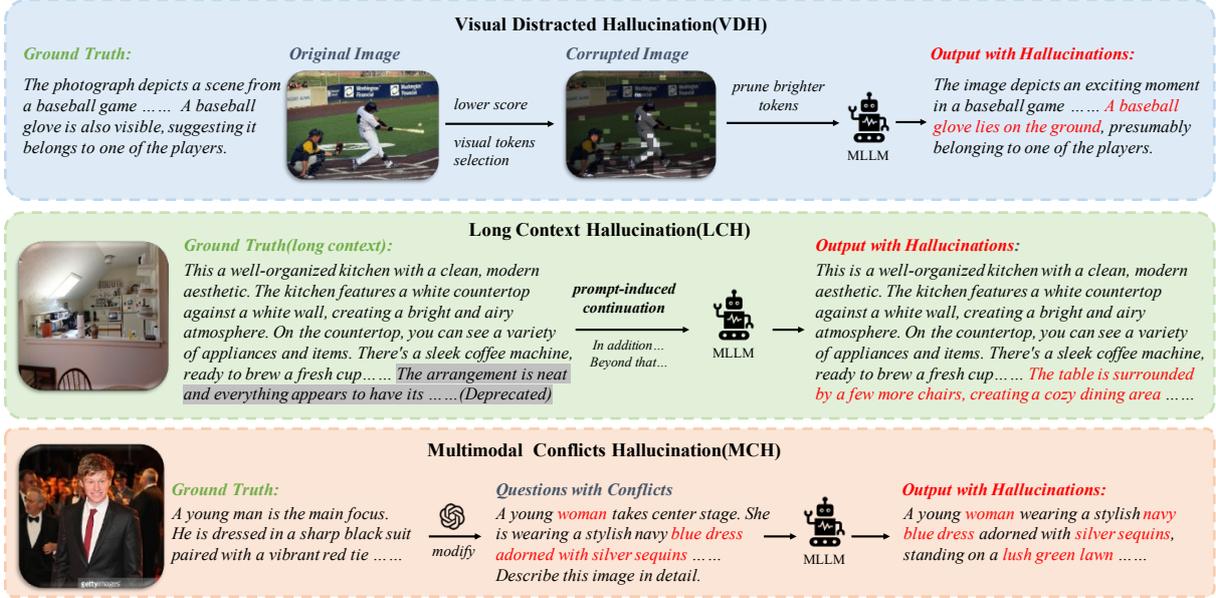


Figure 1: Overview of our three kinds of Hallucinated-targeted Preference data. Better view on the digital screen.

To this end, we develop three types of pairwise samples specifically targeting hallucination issues: Visual Distracted Hallucination (VDH), Long Context Hallucination (LCH), and Multimodal Conflict Hallucination (MCH). An overview of each data type is provided in fig. 1, and further details are outlined in the sections below.

3.3 Visual Distracted Hallucination

Previous works generate negative samples by adding noise to create blurred images, while it may not always produce sufficiently effective negative samples, as indicated in appendix B. A more straightforward way is to construct negative samples using prompts, but the negative samples generated under prompt interference may fail to accurately reflect the issues related to the visual capabilities of MLLMs.

Therefore, to more precisely capture the insufficient visual capabilities of MLLMs, we propose more carefully designed novel approaches from attention perspective. Inspired by SID, we induce vision-and-text association hallucinations by leveraging vision tokens with low attention scores in the self-attention module. Formally, for the transformer block in the auto-regressive decoder, text instructions, vision inputs, and generated tokens are concatenated and projected into three vectors: Q, K and V. The self-attention computes the relevance of each element to the others as follows to get the attention matrix:

$$\mathbf{A} = \text{softmax}((\mathbf{Q} \cdot \mathbf{K}^T) / \sqrt{d} + \mathbf{M}) \quad (1)$$

where d represents the dimension of \mathbf{Q} , \mathbf{K} , \mathbf{V} , \mathbf{M} represents the casual mask. $\mathbf{A} \in R^{(b,h,n,n)}$, where b , h , and n denote batch size, number of key-value heads, and total token number, respectively. We denote the \mathbf{A}_i as the attention matrix after Layer i of MLLMs. Then we calculate vision token importance scores ($\text{Score}_i(v)$) based on \mathbf{A}_i :

$$\text{Score}_i(v) = \frac{1}{h} \sum_{j=1}^h \mathbf{A}_i^{(:,j,:)}[-1] \quad (2)$$

During the model’s auto-regressive decoding process, we retain the K vision tokens with the lowest importance scores, and the resulting decoded response serves as negative samples. By removing the most important visual token from the model in this way, amplifies the influence of relatively irrelevant visual tokens, thus constructing visual information distracted hallucinations as negative samples, urging MLLMs to pay attention to more important visual information.

3.4 Long Context Hallucination

As previously mentioned, the occurrence of hallucinations tends to increase as models generate longer responses. To illustrate this more clearly, we present CHAIR scores by varying the ‘max new tokens’ parameter. As shown in fig. 2, the CHAIR score of LLaVA-v1.5-7B exhibits a clear positive correlation with the ‘max new tokens’, indicating that more hallucinations are produced as the generated content increases. This issue has also been highlighted in recent studies (Yue et al., 2024).

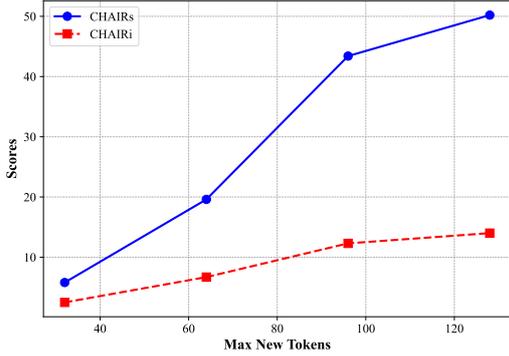


Figure 2: CHAIR scores under different max new tokens

This phenomenon is both logical and explainable. As the model generates longer texts, the proportion of text tokens gradually increases while the proportion of image tokens decreases. This shift causes the model to increasingly neglect visual tokens, resulting in descriptions that appear reasonable but fail to accurately align with the visual content. Our aim is to construct preference data that guides the model to better align its generated content with the visual input and the given question, even when generating long responses. However, existing datasets lack sufficient positive and negative pairs for long-form content and often contain noise with other factors, making them difficult to directly apply for training. *To address this, we firstly propose approach for constructing positive and negative preference pairs for long-form content, ensuring the long text hallucinations while maintaining minimal semantic divergence.*

Given our focus on relatively long-form content, the responses need to be sufficiently lengthy (high-quality long responses). For negative examples, we truncate the last two sentences from a positive example and use the preceding portion as a prefix. The MLLM then continues generating text from this prefix, which compels the model to produce common errors associated with extended text generation. This process is repeated by concatenating the newly generated content to the prefix for three iterations in a loop.

Hint Phrase. Simply providing the prefix and instructing the model to continue often results in unexpected behavior, as the model tends to conclude the response quickly, generating low-information descriptions. To address this issue, we append a 'hint phrase' to the prefix, guiding the model toward producing more informative and detailed responses. Besides that, we also modify the system prompt. Details can be seen in appendix D.2. It helps produce responses prone to more likely errors when

generating long texts. By creating positive and negative pairs in this manner, we aim to use DPO to teach the model how to minimize hallucinations in long-form responses and improve alignment.

3.5 Multimodal Conflicts Hallucination

One of the more challenging yet often overlooked scenarios in mainstream evaluation tasks involves conflicts between modalities. In such cases, models may naturally favor textual content due to their autoregressive generating manner and the larger proportion of the language model component, leading to incorrect outputs. *In this paper, we bring this issue to the forefront to address and firstly use preference optimization to mitigate it.*

To be specific, we construct positive and negative pairs with conflicting prefixes and apply DPO to optimize the model. Specifically, we utilize GPT-4o-mini to rewrite details of the positive examples through prompting, generating information conflicting with the image contents. These conflicting informations are then placed at the beginning of normal questions, prompting the model to produce incorrect responses. As shown in fig. 3, the model is indeed prone to being hallucinated by the conflicting prefixes. We take the model's incorrect outputs as negative examples. Further details on the prompts can be found in fig. 9. Unlike previous types of data, the questions for training of MCH contain conflicting prefixes, as we aim for the model to generate correct responses in the query even when presented with conflicting information.

3.6 Implement details

For LCH, which requires longer responses, we sampled 6k examples with over 300 tokens from ShareGPT4V. For MCH, we randomly sampled 6k examples from ShareGPT4V. For VDH, we obtain 6k examples from ShareGPT4V and 4k examples from VG with positive examples from HA-DPO to enhance data diversity; the preserved K is 500, with other settings aligned with SID (e.g., $i = 2$). Details of data can be found in appendix D.

4 Experiments

In this section, we empirically investigate the evaluation of HDPO. We begin by describing the experimental settings, including the evaluation datasets and training details. Next, we present the results on various hallucination evaluation datasets, demonstrating the promising performance of HDPO. Additionally, we validate the expected functions of

	POPE	CHAIR		AMBER				
	F1 Score \uparrow	CHAIR _s \downarrow	CHAIR _i \downarrow	CHAIR \downarrow	HalRate \downarrow	Cog. \downarrow	F1 Score \uparrow	AMBER-S \uparrow
LLaVA-v1.5-7B	86.1	51.2	14.2	7.6	35.1	4.3	74.5	83.5
Vlfeedback \dagger	83.7	40.3	13.2	–	–	–	–	–
POVID \dagger	86.9	35.2	8.3	–	–	–	–	–
HA-DPO	86.9	37.2	10.0	6.4	29.9	3.2	78.2	85.9
SeVa	86.8	54.6	15.9	7.4	35.6	3.2	84.1	88.3
BPO	83.1	42.2	10.1	5.0	33.5	2.0	84.5	89.7
CSR	87.0	19.6	5.4	3.8	16.9	1.4	76.0	86.1
HDPO (ours)	86.8	16.6	5.1	3.3	15.8	0.8	84.1	90.4

Table 1: Experimental results of HDPO on LLaVA-v1.5-7B compared with baselines applied on LLaVA-v1.5-7B. The best result for each metric is in bold. Some results \dagger are referenced from Zhou et al. (2024b). The F1 of POPE and AMBER are discriminative metrics, AMBER-s is a comprehensive metric, and the others are generative metrics.

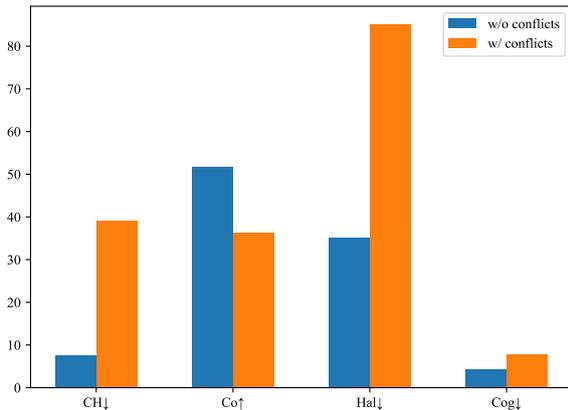


Figure 3: Performance of LLaVA-v1.5-7B w/ and w/o conflicts on AMBER, details in section 4.4.2.

LCH and MCH. Finally, we provide ablation studies and conduct in-depth analyses in more detail.

4.1 Experimental Settings

Evaluation Datasets. We evaluate the effectiveness of HDPO in mitigating hallucinations across both captioning tasks and simplified visual question answering (VQA) tasks using three evaluation datasets as follows: (1) CHAIR is an evaluation method used in image captioning tasks to assess object hallucinations in model responses. There are two metrics: CHAIR_s and CHAIR_i. CHAIR_s measures hallucinations at the sentence level, while CHAIR_i measures them at the image level respectively. (2) POPE is a popular dataset for evaluating object hallucinations in MLLMs. We calculate and report the average F1 score on different splits. (3) AMBER is an LLM-free multi-dimensional benchmark, offering a cost-effective and efficient evaluation pipeline. It supports the evaluation of both generative and discriminative tasks including hallucinations related to existence, attributes, and relations. For all details of datasets and metrics used can be seen in appendix A.

Training Details. As most related works (Chen

et al., 2023; Zhu et al., 2024; Pi et al., 2025) are carried on LLaVA-v1.5 (Liu et al., 2024a), we select it as our base model for experiments, which allows for easy comparison with other existing works. Models’ weights are pretrained and further fine-tuned using supervised fine-tuning (SFT) before applying HDPO. During the training phase, we employ Zero stage-3 optimization and use Vicuna-7B/13B and CLIP-VIT-L-336px as our LLM and vision encoder, respectively. The training is conducted with 2 epochs with a batch size of 64, a learning rate of 2e-6, weight decay as 0, LoRA rank as 64, and a beta value of 0.1. All experiments are run on one single machine with 8 A800 GPUs. The total training time is 3 hours for LLaVA-v1.5-7B and 4 hours for LLaVA-v1.5-13B. Besides, we also validate HDPO on InstructBLIP, further demonstrating effectiveness in section 4.3.

Competitor. We first compare HDPO with its base model. We also select several preference learning methods, including Vlfeedback (Li et al., 2024), POVID (Zhou et al., 2024a), CLIP-DPO (Ouali et al., 2025), HA-DPO (Zhao et al., 2023), SeVa (Zhu et al., 2024), BPO (Pi et al., 2025), and CSR (Zhou et al., 2024b). Furthermore, we compare HDPO on AMBER with other MLLMs in appendix D.5.

4.2 Results on Diverse Hallucination Tasks

HDPO achieves SOTA level on both generative and discriminative hallucination tasks. The results indicate that HDPO performs well in mitigating hallucinations, achieving almost SOTA level, especially on generative tasks. This outcome is natural, as our data contains only descriptive content, leading to relatively strong performance on generative tasks. Since we don’t specifically construct data tailored for discriminative tasks, the improvement in these tasks is not substantial. However,

	POPE	CHAIR		AMBER				
	F1 Score \uparrow	CHAIR _s \downarrow	CHAIR _i \downarrow	CHAIR \downarrow	HalRate \downarrow	Cog. \downarrow	F1 Score \uparrow	AMBER-S \uparrow
LLaVA-v1.5-13B	85.8	48.0	13.6	6.6	31.0	3.3	73.0	83.2
HA-DPO	87.3	46.0	12.1	6.0	30.7	3.0	79.1	86.6
SeVa	86.9	59.8	17.4	9.0	43.3	3.7	84.8	87.9
CSR	87.3	24.0	5.6	3.6	19.0	1.8	73.1	84.8
HDPO (ours)	87.6	15.4	5.3	3.8	16.5	0.8	81.2	88.7

Table 2: Experimental results of HDPO on LLaVA-v1.5-13B compared with baselines applied on LLaVA-v1.5-13B. More details of baselines can be seen in appendix C.

the overall performance remains strong, indicating that our approach, which targets the sources of hallucinations rather than specific tasks, is more effective for mitigating hallucinations. Notably, HDPO achieves **67.6%** improvement on CHAIR_s, **64.1%** improvement on CHAIR_i, **55%** enhancement on HalRate, best performance on AMBER-S. Besides, we also evaluate HDPO on a comprehensive benchmark, MM-Vet (Yu et al., 2024), where we observe a slight improvement. This aligns with our expectations, as the model is not fine-tuned on a wide range of tasks and data types, but focused on reducing hallucinations.

Brief analyses on other baselines. Some baselines lack comprehensive performance on hallucination evaluation. SeVa, though effective on AMBER’s discriminative tasks, shows no improvement on generative tasks, likely due to its reliance on VQA-type data. Similarly, BPO underperforms on CHAIR. In contrast, CSR excels in generative tasks but struggles with AMBER’s discriminative tasks. This indicates that while these methods enhance model performance, they do not fully optimize for hallucination, and their ability to mitigate hallucinations remains inconsistent and incomplete, while HDPO demonstrates strong performance in hallucination evaluation, as evidence of its ‘hallucination-targeted’ design.

Advantages of our HDPO Data. The size of our dataset also provides a relative advantage. For instance, with nearly 12% data amount compared with BPO, HDPO significantly improves model’s performance on hallucination, achieving better performance than BPO on generative tasks by a large margin. Moreover, we did not construct VQA data for discriminative tasks. Nevertheless, the results are already impressive, demonstrating that our HDPO is universally effective.

4.3 Universality on Different Base Models

We also conduct experiments across different base models to verify our HDPO’s universality. Specifi-

cally, we apply HDPO to the widely-used LLaVA-v1.5-13B for MLLM hallucination evaluation. The results are shown in table 2, demonstrating that the model’s performance remains consistent with expectations, with improvements in hallucination mitigation. It also implies that our generated hallucination-targeted DPO data is effective for different LLM sizes. To further validate the generalization capabilities of other MLLMs, we also conduct experiments on InstructBLIP (Liu et al., 2024b). The results in table 5 also show consistent improvement on the overall performance.

4.4 Analyses on Different Hallucinations

The results from above experiments demonstrate our method’s superior performance in mitigating hallucinations. However, do they truly work effectively in the scenarios we claim? Below, we briefly design two more challenging sub-tasks of hallucination that align with our claims, aiming to further showcase the effectiveness of our data construction of LCH and MCH. We also conduct experiments to compare VDH with adding noise in appendix B, further demonstrating effectiveness of VDH.

4.4.1 Long Context Hallucination

To evaluate the effectiveness of LCH on longer responses, we conduct an extended experiment on the AMBER generative task. Specifically, when the model is asked the question "Describe this image in detail", we append the instruction "answer in 800 words" to encourage longer responses. As indicated in table 3, HDPO shows good and stable performance in handling longer responses, with the lowest HalRate, CHAIR_s, and Cog. It demonstrates that our construction for LCH works as expected in longer responses.

4.4.2 Multimodal Conflicts Hallucination

In real-world scenarios, multimodal conflicts are common when using MLLMs. To better evaluate the model’s performance under such conditions, we

	CHAIR ↓	HalRate ↓	Cog. ↓
LLaVA-v1.5-7B	9.0	45.1	5.7
HA-DPO	7.5	37.6	4.4
SeVa	7.5	43.4	4.3
BPO	6.4	55.3	4.8
HDPO	3.4	21.4	1.3
w/o LCH	4.6	26.4	1.8

Table 3: Results of long context hallucination.

	CHAIR ↓	HalRate ↓	Cog. ↓
LLaVA-v1.5-7B	39.1	85.1	7.8
HA-DPO	40.3	86.1	8.1
SeVa	39.1	86.1	7.8
BPO	22.3	81.2	7.7
HDPO	14.3	52.0	5.2
w/o MCH	39.8	84.7	6.7

Table 4: Results of multimodal conflict hallucination.

design a more challenging task. Specifically, we randomly select 200 questions from the generative task in the AMBER dataset. First, LLaVA-1.5-7B is used to generate answers for these questions to get coarse-grained image descriptions. Next, GPT-4o-mini rewrites the details in the descriptions, following the construction method of MCH. We then introduce the incorrect information as a prefix to the question and ask the model to describe the image while influenced by the conflicting context.

The experimental results are shown in table 4, demonstrating that despite encountering conflicting prefixes, our HDPO maintains promising performance. Compared to other baselines, HDPO achieves the best scores in CHAIR_s, HalRate, and Cog. It reveals that our HDPO shows significant improvement in the model’s performance under this more difficult setting, highlighting the effectiveness of MCH. Additionally, we also make a comparison between the effects of adding noise and preserved visual tokens with lower scores. Further details can be seen in the appendix B.

4.5 Ablation Study

To demonstrate the contributions of VDH, LCH, and MCH to overall performance, we progressively remove each component and report the results. (1) As shown in table 6, the performance declines as we remove each data type. The model achieves the best performance when all three data types are included. These experimental results confirm the individual contributions of each component. (2) It can also be observed that after incorporating MCH, there is no improvement in CHAIR_s and CHAIR_i. However, the inclusion of both posi-

	POPE ↑	CHAIR _s ↓	CHAIR _i ↓	AMBER-S ↑
InstructBLIP	83.7	57.0	16.1	82.5
HA-DPO	85.6	56.6	15.5	84.3
HDPO (ours)	84.8	34.8	10.9	85.9

Table 5: Results of HDPO on InstructBLIP-13B.

	CHAIR		AMBER	
	CHAIR _s ↓	CHAIR _i ↓	CHAIR ↓	F1 ↑
LLaVA-v1.5-7B	51.4	14.2	7.6	74.5
+VDH +LCH +MCH	16.6	5.1	3.3	84.1
+LCH +MCH	28.4	7.5	4.8	78.9
+MCH	51.2	15.1	7.6	78.1

Table 6: Results of ablation study.

tive and negative examples for training leads to improvement in F1 of discriminative task (**4.8%**↑). (3) With the addition of LCH, F1 of the discriminative task shows minimal change, whereas the generative task demonstrates a substantial improvement, with CHAIR_s (**44.5%**↓) and CHAIR_i (**50.3%**↓) showing marked gains. This indicates that LCH is particularly effective for generative tasks. (4) Finally, incorporating VDH enhances model’s performance across all tasks, and the combination of all three categories achieves the best results. The significance of LCH and MCH is also verified in section 4.4 with the corresponding tasks.

4.6 Scalability of HDPO

We analyze the impact of data size on our method. The performance of LLaVA-v1.5-7B fine-tuned on datasets of varying sizes with the same proportions are shown in fig. 4. As the data size increases, the effectiveness of our approach also improves, highlighting the potential for scaling up. This demonstrates the superior performance of HDPO.

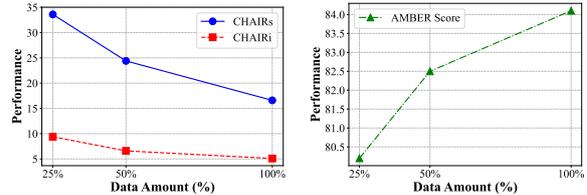


Figure 4: Scalability of HDPO with different data sizes.

5 Conclusion

In this paper, we present HDPO, a novel approach designed to effectively mitigate hallucinations in MLLMs. We analyze three types of hallucinations observed in MLLMs and create hallucination preference data based on the identified causes. Extensive experiments across different benchmarks demonstrate the ability of HDPO to reduce hallucinations in MLLMs, showing effectiveness.

634 Limitations

635 In this paper, we introduce HDPO, which effectively mitigates the hallucination problem in current multimodal large language models. However, several issues remain unresolved. Specifically, we have not yet developed distinct strategies for controlling data quality, and the generation of automated negative examples lacks methods for further verification and optimization, which could improve the effectiveness of our approach. Additionally, there may be opportunities to further enhance the quality of positive examples. Moreover, our construction methods and strategies could potentially be integrated with other techniques for processing more high-quality preference data, which may further improve the model’s performance. Fine-tuning larger models with extensive, integrated datasets may not only enhance overall reasoning capabilities but also increase the model’s robustness against hallucinations. This represents a promising area for further investigation, and we leave these open questions for future research.

656 Ethics Statement

657 This work mitigates hallucinations of multimodal large language models to enhance their reliability and practicality. We have carefully considered the ethical implications of our work. The models and datasets we used are publicly available and commonly used, and our findings may inherit the biases and limitations carried out in these resources.

664 References

665 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

670 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*.

675 Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou. 2024. Hallucination of multimodal large language models: A survey. *arXiv preprint arXiv:2404.18930*.

679 Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023. Sharegpt4v: Improving large multimodal models with better captions. *arXiv preprint arXiv:2311.12793*.

Zhaorun Chen, Zhuokai Zhao, Hongyin Luo, Huaxiu Yao, Bo Li, and Jiawei Zhou. 2024. HALC: Object hallucination reduction via adaptive focal-contrast decoding. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 7824–7846. PMLR. 684 685 686 687 688 689 690

Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zichong Yang, Kuei-Da Liao, et al. 2024. A survey on multimodal large language models for autonomous driving. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 958–979. 691 692 693 694 695 696 697

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*. 698 699 700 701 702

Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. 2024. Opera: Alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13418–13427. 703 704 705 706 707 708 709 710

Fushuo Huo, Wenchao Xu, Zhong Zhang, Haozhao Wang, Zhicheng Chen, and Peilin Zhao. 2024. Self-introspective decoding: Alleviating hallucinations for large vision-language models. *arXiv preprint arXiv:2408.02032*. 711 712 713 714 715

Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. 2024. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13872–13882. 716 717 718 719 720 721 722

Lei Li, Zhihui Xie, Mukai Li, Shunian Chen, Peiyi Wang, Liang Chen, Yazheng Yang, Benyou Wang, Lingpeng Kong, and Qi Liu. 2024. VLFeedback: A large-scale AI feedback dataset for large vision-language models alignment. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6227–6246, Miami, Florida, USA. Association for Computational Linguistics. 723 724 725 726 727 728 729 730 731

Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Xin Zhao, and Ji-Rong Wen. 2023. 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. 732 733 734 735 736 737

Fenglin Liu, Tingting Zhu, Xian Wu, Bang Yang, Chenyu You, Chenyang Wang, Lei Lu, Zhangdaihong Liu, Yefeng Zheng, Xu Sun, et al. 2023a. A 738 739 740

741	medical multimodal large language model for future	Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma,	796
742	pandemics. <i>NPJ Digital Medicine</i> , 6(1):226.	Yann LeCun, and Saining Xie. 2024. Eyes wide	797
743	Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser	shut? exploring the visual shortcomings of multi-	798
744	Yacoub, and Lijuan Wang. 2023b. Mitigating hal-	modal llms. In <i>Proceedings of the IEEE/CVF Con-</i>	799
745	lucination in large multi-modal models via robust	<i>ference on Computer Vision and Pattern Recognition</i> ,	800
746	instruction tuning. In <i>The Twelfth International Con-</i>	pages 9568–9578.	801
747	<i>ference on Learning Representations</i> .	Junyang Wang, Yuhang Wang, Guohai Xu, Jing Zhang,	802
748	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	Yukai Gu, Haitao Jia, Ming Yan, Ji Zhang, and Ji-	803
749	Lee. 2024a. Visual instruction tuning. <i>Advances in</i>	tao Sang. 2023a. An llm-free multi-dimensional	804
750	<i>neural information processing systems</i> , 36.	benchmark for mllms hallucination evaluation. <i>arXiv</i>	805
751	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	<i>preprint arXiv:2311.07397</i> .	806
752	Lee. 2024b. Visual instruction tuning. <i>Advances in</i>	Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi	807
753	<i>neural information processing systems</i> , 36.	Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,	808
754	Jiazhen Liu, Yuhan Fu, Ruobing Xie, Runquan Xie,	Lei Zhao, Xixuan Song, et al. 2023b. Cogvlm: Vi-	809
755	Xingwu Sun, Fengzong Lian, Zhanhui Kang, and	sual expert for pretrained language models. <i>arXiv</i>	810
756	Xirong Li. 2024c. Phd: A prompted visual	<i>preprint arXiv:2311.03079</i> .	811
757	hallucination evaluation dataset. <i>arXiv preprint</i>	Junfei Wu, Qiang Liu, Ding Wang, Jinghao Zhang, Shu	812
758	<i>arXiv:2403.11116</i> .	Wu, Liang Wang, and Tienyi Tan. 2024. Logical	813
759	OpenAI. 2023. GPT-4V(ision) system card .	closed loop: Uncovering object hallucinations in	814
760	OpenAI. 2024. Hello GPT-4o .	large vision-language models . In <i>Findings of the As-</i>	815
761	Yassine Ouali, Adrian Bulat, Brais Martinez, and	<i>sociation for Computational Linguistics: ACL 2024</i> ,	816
762	Georgios Tzimiropoulos. 2025. Clip-dpo: Vision-	pages 6944–6962, Bangkok, Thailand. Association	817
763	language models as a source of preference for fix-	for Computational Linguistics.	818
764	ing hallucinations in lvlms. In <i>Computer Vision –</i>	Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, An-	819
765	<i>ECCV 2024</i> , pages 395–413, Cham. Springer Nature	wen Hu, Haowei Liu, Qi Qian, Ji Zhang, and Fei	820
766	Switzerland.	Huang. 2024. mplug-owl2: Revolutionizing multi-	821
767	Renjie Pi, Tianyang Han, Wei Xiong, Jipeng Zhang,	modal large language model with modality collabora-	822
768	Runtao Liu, Rui Pan, and Tong Zhang. 2025.	tion. In <i>Proceedings of the IEEE/CVF Conference</i>	823
769	Strengthening multimodal large language model with	<i>on Computer Vision and Pattern Recognition</i> , pages	824
770	bootstrapped preference optimization. In <i>European</i>	13040–13051.	825
771	<i>Conference on Computer Vision</i> , pages 382–398.	Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao	826
772	Springer.	Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun,	827
773	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	and Enhong Chen. 2023. Woodpecker: Hallucina-	828
774	pher D Manning, Stefano Ermon, and Chelsea Finn.	tion correction for multimodal large language models.	829
775	2024. Direct preference optimization: Your language	<i>arXiv preprint arXiv:2310.16045</i> .	830
776	model is secretly a reward model. <i>Advances in Neu-</i>	Weihaio Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang,	831
777	<i>ral Information Processing Systems</i> , 36.	Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan	832
778	Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns,	Wang. 2024. MM-vet: Evaluating large multimodal	833
779	Trevor Darrell, and Kate Saenko. 2018. Object hallu-	models for integrated capabilities . In <i>Proceedings of</i>	834
780	ciation in image captioning . In <i>Proceedings of the</i>	<i>the 41st International Conference on Machine Learn-</i>	835
781	<i>2018 Conference on Empirical Methods in Natural</i>	<i>ing</i> , volume 235 of <i>Proceedings of Machine Learning</i>	836
782	<i>Language Processing</i> , pages 4035–4045, Brussels,	<i>Research</i> , pages 57730–57754. PMLR.	837
783	Belgium. Association for Computational Linguistics.	Zihao Yue, Liang Zhang, and Qin Jin. 2024. Less is	838
784	Min Shi, Fuxiao Liu, Shihao Wang, Shijia Liao, Sub-	more: Mitigating multimodal hallucination from an	839
785	hashree Radhakrishnan, De-An Huang, Hongxu Yin,	EOS decision perspective . In <i>Proceedings of the</i>	840
786	Karan Sapra, Yaser Yacoub, Humphrey Shi, et al.	<i>62nd Annual Meeting of the Association for Computa-</i>	841
787	2024. Eagle: Exploring the design space for multi-	<i>tional Linguistics (Volume 1: Long Papers)</i> , pages	842
788	modal llms with mixture of encoders. <i>arXiv preprint</i>	11766–11781, Bangkok, Thailand. Association for	843
789	<i>arXiv:2408.15998</i> .	Computational Linguistics.	844
790	Xingwu Sun, Yanfeng Chen, Yiqing Huang, Ruobing	Zhiyuan Zhao, Bin Wang, Linke Ouyang, Xiaoyi Dong,	845
791	Xie, Jiaqi Zhu, Kai Zhang, Shuai-peng Li, Zhen Yang,	Jiaqi Wang, and Conghui He. 2023. Beyond hallu-	846
792	Jonny Han, Xiaobo Shu, et al. 2024. Hunyuan-	ciations: Enhancing lvlms through hallucination-	847
793	large: An open-source moe model with 52 billion	aware direct preference optimization. <i>arXiv preprint</i>	848
794	activated parameters by tencent. <i>arXiv preprint</i>	<i>arXiv:2311.16839</i> .	849
795	<i>arXiv:2411.02265</i> .	Yiyang Zhou, Chenhang Cui, Rafael Rafailov, Chelsea	850
		Finn, and Huaxiu Yao. 2024a. Aligning modalities	851
		in vision large language models via preference fine-	852
		tuning. <i>arXiv preprint arXiv:2402.11411</i> .	853

854 Yiyang Zhou, Zhiyuan Fan, Dongjie Cheng, Sihan Yang,
855 Zhaorun Chen, Chenhang Cui, Xiyao Wang, Yun
856 Li, Linjun Zhang, and Huaxiu Yao. 2024b. [Cali-](#)
857 [brated self-rewarding vision language models](#). In
858 *The Thirty-eighth Annual Conference on Neural In-*
859 *formation Processing Systems*.

860 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and
861 Mohamed Elhoseiny. 2023. Minigt-4: Enhancing
862 vision-language understanding with advanced large
863 language models. *arXiv preprint arXiv:2304.10592*.

864 Ke Zhu, Liang Zhao, Zheng Ge, and Xiangyu Zhang.
865 2024. [Self-supervised visual preference alignment](#).
866 In *Proceedings of the 32nd ACM International Con-*
867 *ference on Multimedia, MM '24*, page 291–300, New
868 York, NY, USA. Association for Computing Machin-
869 ery.

870 Zhuofan Zong, Bingqi Ma, Dazhong Shen, Guanglu
871 Song, Hao Shao, Dongzhi Jiang, Hongsheng Li,
872 and Yu Liu. 2024. Mova: Adapting mixture of vi-
873 sion experts to multimodal context. *arXiv preprint*
874 *arXiv:2404.13046*.

A Details of Datasets and Metrics 875

876 We evaluate the effectiveness of HDPO in mitigat-
877 ing hallucinations across both captioning tasks and
878 simplified visual question answering (VQA) tasks
879 using three evaluation datasets as follows:

880 **CHAIR** (Rohrbach et al., 2018): The Caption
881 Hallucination Assessment with Image Relevance
882 (CHAIR) is an evaluation method used in image
883 captioning tasks to assess object hallucinations in
884 model responses. There are two metrics: CHAIR_s
885 and CHAIR_i. CHAIR_s measures hallucinations at
886 the sentence level, while CHAIR_i measures them
887 at the image level respectively. We conduct the
888 CHAIR evaluation on the MSCOCO dataset follow-
889 ing the setting in OPERA(Huang et al., 2024) with
890 500 random images. For each image, the model is
891 prompted with: "Please describe this image in de-
892 tail." to obtain their descriptions. By default, we set
893 the 'max new tokens' to 512. More specifically, the
894 calculation for the CHAIR_s and CHAIR_i metrics
895 are as follows:

$$\text{CHAIR}_s = \frac{|\{\text{hallucinated objects}\}|}{|\{\text{all mentioned objects}\}|} \quad (3) \quad 896$$

$$\text{CHAIR}_i = \frac{|\{\text{captions w/ hallucinated objects}\}|}{|\{\text{all captions}\}|} \quad (4) \quad 897$$

898 **POPE** (Li et al., 2023): The Polling-based Object
899 Probing Evaluation (POPE) is a popular dataset for
900 evaluating object hallucinations in MLLMs. The
901 evaluation is asking the model questions in the
902 format: "Is there a <object> in the image?". It can
903 be divided into three splits: popular, adversarial,
904 and random. In the popular split, the evaluation
905 targets the most frequently occurring objects in
906 the dataset. In the adversarial split, it assesses the
907 MLLM's ability to identify objects that are highly
908 relevant to those present in the image. We evaluate
909 the metrics for all splits, and calculate and report
910 the average F1 score. POPE can be constructed on
911 different datasets, and we evaluate models on the
912 POPE dataset built on COCO. 913

914 **AMBER** (Wang et al., 2023a): An Automated
915 Multi-dimensional Benchmark for Multi-modal
916 Hallucination Evaluation (AMBER) is an LLM-
917 free multi-dimensional benchmark, offering a cost-
918 effective and efficient evaluation pipeline. It sup-
919 ports the evaluation of both generative and discrim-
920 inative tasks including hallucinations related to ex-
921 istence, attributes, and relations. Its generative eval-
922 uation aligns with our desired assessment of long

descriptions, while the other dimensions provide insights into the model’s performance on relatively simple VQA tasks, thereby reflecting the model’s hallucination comprehensively. For its generative task, three metrics are used: CHAIR, Hal, and Cog. **CHAIR** measures the frequency of hallucinatory objects in the responses, **Hal** represents the proportion of responses containing hallucinations, and **Cog** assesses whether the hallucinations produced by MLLMs resemble those found in human cognition. For its discriminative task, we calculate and report the average F1 score. We also calculate **AMBER Score** denoted as AMBER-S, which reflects overall performance, and it’s calculated as follows:

$$AMBER\ Score = \frac{1}{2} \times (1 - CHAIR + F1) \quad (5)$$

B Comparison of noise and token preservation

We also conduct experiments to compare the impact of adding noise versus preserving visual tokens. Specifically, we use 6k samples from ShareGPT4V to construct negative samples by introducing diffusion noise and preserving visual tokens, and train the LLaVA-v1.5-7B model by direct preference optimization. The results of these experiments are presented in table 7. As the experimental results show, using visual token preservation can achieve better performance on hallucination evaluation.

C Baseline Selection of 13B

For the experiments on the 13B model, we select several recent strong baselines, including SeVa and CSR, using their open-sourced checkpoints for evaluation. Additionally, we reimplement HA-DPO on LLaVA-v1.5-13B, as the original repository does not provide this checkpoint. We also attempt to reimplement BPO on LLaVA-v1.5-13B with no available checkpoints, the evaluation results are unexpectedly low, with POPE scores falling below 80.0. Therefore, these results are not included in the table. However, the BPO results for the 7B model are obtained using the publicly released checkpoints. For InstructBLIP, we don’t find other preference optimization works on it.

D Details about Our data

D.1 Visual Distracted Hallucination

We obtain positive examples for our dataset from two sources: VG(with positive examples in HA-

DPO) and ShareGPT4V. After extracting positive examples from ShareGPT4V, we found them to be too long. To mitigate length bias, we used GPT4o-mini to rewrite them to match the length of the negative examples. The prompt used is shown in fig. 7. For positive examples sourced from HA-DPO, after generating negative examples, we followed the original approach by rewriting the negative examples using GPT4o-mini. The prompt used is shown in fig. 6. Also, we can adopt the method in HA-DPO to create more data. For k and i , we make an empirical choice based on performance and original settings.

D.2 Long Context Hallucination

We use LLaVA-1.5-7B to continue generating text for the positive examples, with the system prompt in fig. 5, and the hint phrases in fig. 8. By excluding the last two sentences, we aim to increase the concentration of hallucinated content in the tail of the response. Generating three continuations at a time maintains an approximate balance in the average length between positive and negative examples.

D.3 Multimodal Conflicts Hallucination

We utilize GPT-4o-mini to modify the details of the positive examples, following the prompt shown in fig. 9. This approach introduces conflicting information that deviates from the image content.

D.4 Effect of data ratio

We did not conduct detailed experiments comparing different data type ratios. However, throughout the experiments, all tested ratios showed significant improvements over the original model. We report the best-performing dataset from our experiments. Determining the optimal ratio of different data types is inherently a more challenging and general problem, which goes beyond the scope of this paper.

D.5 Comparison on AMBER with other MLLMs

We also report the hallucination evaluation results on AMBER for both generative and discriminative tasks of HDPO on LLaVA-1.5-7B compared with other MLLMs including mPLUG-Owl2 (Ye et al., 2024), MiniGPT4 (Zhu et al., 2023), CogVLM (Wang et al., 2023b), Qwen-VL (Bai et al., 2023) and GPT4V (OpenAI, 2023) in table 9.

	POPE	CHAIR		AMBER				
	F1 Score \uparrow	CHAIR _s \downarrow	CHAIR _i \downarrow	CHAIR \downarrow	HalRate \downarrow	Cog. \downarrow	F1 Score \uparrow	AMBER-S \uparrow
LLaVA-v1.5-7B	86.1	51.2	14.2	7.6	35.1	4.3	74.5	83.5
+ Diffu _{6k}	86.2	62.8	18.4	9.2	47.5	4.3	78.1	84.5
+ VDH _{6k}	87.1	48.2	13.7	6.1	32.0	2.7	80.2	87.1

Table 7: Experimental results of LLaVA-v1.5-7B trained with two ways to construct preference pairs: adding noise and preserving visual tokens. The diffusion noise step is 800. The best result for each metric is in bold.

	Len	Cover	Co. / Len \uparrow	CHAIR \downarrow
LLaVA-1.5-7B	75.0	51.8	0.69	7.6
BPO	148.0	58.8	0.40	5.0
SeVa	76.0	53.4	0.70	7.4
CSR	64.0	45.0	0.70	3.8
HDPO	69.0	50.2	0.73	3.3

Table 8: Analysis of Cover. on AMBER

D.6 Computational cost and efficiency Compared with Baselines

As computational efficiency is critical for real-world applications, we present the training costs of HDPO and other baseline methods as follows.

CSR: Training utilized one A100 GPU, with LLaVA-1.5 (7B / 13B) fine-tuned for approximately 3.5 / 5.0 hours.

SeVa: Training utilized 8 A800 GPUs, with LLaVA-1.5 (7B / 13B) fine-tuned for approximately 0.7 / 1.3 hours.

BPO: Training utilized 8 A40 GPUs, with LLaVA-1.5 (7B / 13B) fine-tuned for approximately 17.0 / 28.0 hours.

HDPO: Training utilized 8 A800 GPUs, with LLaVA-1.5 (7B / 13B) fine-tuned for approximately 3.0 / 4.0 hours.

Training time is fundamentally influenced by the size of the training dataset. Except for BPO, which requires a relatively longer training time, the training costs and durations for the other methods fall within a comparable range. Thus, we believe that our method holds significant value for practical applications.

D.7 More Analysis of Cover

There is another Cover metric in AMBER, represents object coverage. It’s related to the length of generated content. We calculate the Cover / Length and report it in table 8. It shows that HDPO’s outputs are more precise and of higher quality with the highest Co./ Len. Additionally, we have conducted experiments showing that generating longer outputs improves Cover while maintain good hallucination performance.

	CHAIR \downarrow	Hal \downarrow	Cog. \downarrow	F1 \uparrow	AMBER-S \uparrow
mPLUG-Owl	21.6	76.1	11.5	18.9	48.7
LLaVA	11.5	48.8	5.5	32.7	60.6
MiniGPT4	13.6	65.3	11.3	64.7	75.6
CogVLM	5.6	23.6	1.3	72.3	83.4
mPLUG-Owl2	10.6	39.9	4.5	78.5	84.0
Qwen-VL	5.5	23.6	1.9	84.9	89.7
GPT-4V	4.6	30.7	2.6	87.4	91.4
HDPO	3.3	15.8	0.8	84.1	90.4

Table 9: Comparison on AMBER with more MLLMs, most results are source from(Wang et al., 2023a).

D.8 Further Discussion of Limitation

Although HDPO enjoys promising performance in Mitigating Hallucination, there are still some potential boundaries we meet as follows:

(1) For relatively long content generation, HDPO may still struggle to fully address the issue. As the generated content becomes longer, hallucinations may persist. To completely resolve this problem, the model’s intrinsic long-context processing capabilities might first need to be enhanced. However, the current long-text abilities of MLLMs are not as advanced as those of LLMs, which presents an intriguing direction for future exploration.

(2) Additionally, as highlighted by (Tong et al., 2024; Zong et al., 2024; Shi et al., 2024), the visual encoder in current MLLMs operates at a relatively coarse granularity, resulting in insufficient or suboptimal visual features. These limitations cannot be fully addressed by HDPO and will likely require either more powerful visual encoders or improved MLLM architectures, both of which are also promising directions for future research.

(3) What’s more, fine-tuning larger models on extensive, integrated datasets could improve reasoning capabilities and robustness against hallucinations. These open questions remain promising directions for future research.

We think the above additional discussion clarifies the limitations of HDPO and outlines potential directions for addressing these challenges.

System Prompt:

You should describe in detail all elements in the image. Be thorough in addressing aspects such as color, shape, size, position, quantity, actions, emotions, and more. Your response should be as much as possible.

Figure 5: System Prompt used in LCH

Rewrite Prompt:

Help me rewrite the given sentence. Don't change any detail and information in the original sentence. Don't add any new information.

The sentence you need to rewrite: %s
Directly give the rewritten sentence:

Figure 6: Rewrite Prompt used in VDH

Adjust Length Prompt:

Please adjust the length of the Description to approximately %s words.

Ensure all essential details and meanings are preserved, with clear, concise, and accurate expression. Provide the modified Description directly.

Original Description: "%s"
Modified Description:

Figure 7: Adjust Length Prompt used in VDH

Hint Phrases:

"In addition",
"Moreover",
"Furthermore",
"Besides that",
"Additionally",
"What's more",
"As well as that",
"Beyond that",
"There is something else that needs to be mentioned",
"Not only that",
"It should also be noted that"

Figure 8: Hint Phrases used in LCH

Modify Prompt:

I will give you a description of an image, and you need to modify various details of the description, such as the number of objects, types of objects, their positions, colors, behaviors, and so on.

Description: %s
Modified Description:

Figure 9: Modify Prompt used in MCH