# Optimizing LVLMs with On-Policy Data for Effective Hallucination Mitigation

Chengzhi Yu,[1,*]   Yifan Xu,[2,*]   Yifan Chen,[2]   Wenyi Zhang[1,†]

[1] University of Science and Technology of China    [2] Hong Kong Baptist University

## Abstract

*Recently, large vision-language models (LVLMs) have risen to be a promising approach for multimodal tasks. However, principled hallucination mitigation remains a critical challenge. In this work, we first analyze the data generation process in LVLM hallucination mitigation and affirm that on-policy data significantly outperforms off-policy data, which thus calls for efficient and reliable preference annotation of on-policy data. We then point out that, existing annotation methods introduce additional hallucination in training samples, which may enhance the model's hallucination patterns, to address this problem, we propose training a hallucination classifier giving binary annotations, which guarantee clean chosen samples for the subsequent alignment. To further harness of the power of on-policy data, we design a robust iterative direct preference optimization (DPO) algorithm adopting a dynamic sample reweighting scheme. We conduct comprehensive experiments on three benchmarks with comparison to 8 state-of-the-art baselines. In particular, our approach reduces the hallucination rate of LLaVA-1.5-7B on MMHalBench by 50.8% and the average hallucination rate on Object Hal-Bench by 79.5%; more significantly, our method fully taps into the potential of open-source models, enabling LLaVA-1.5-13B to even surpass the performance of GPT-4V.*

## 1. Introduction

The powerful language generation capabilities of pretrained large language models (LLMs) motivates recent advancements of large vision-language models (LVLMs) [1, 5, 20, 21, 38]. However, despite the enhanced capacity empowered by LLMs, current LVLMs are still prone to generating responses that contradict with the reference image [10]. This phenomenon, known as *hallucination*, greatly compromises the generation quality of LVLMs.

In recent studies, *preference alignment* proves to effectively reduce hallucination in LVLMs [17, 48, 49, 53],

---

[*]Equal contribution.

[†]Correspondence to: Wenyi Zhang ⟨wenyizha@ustc.edu.cn⟩.



(a) Prompt and textual context

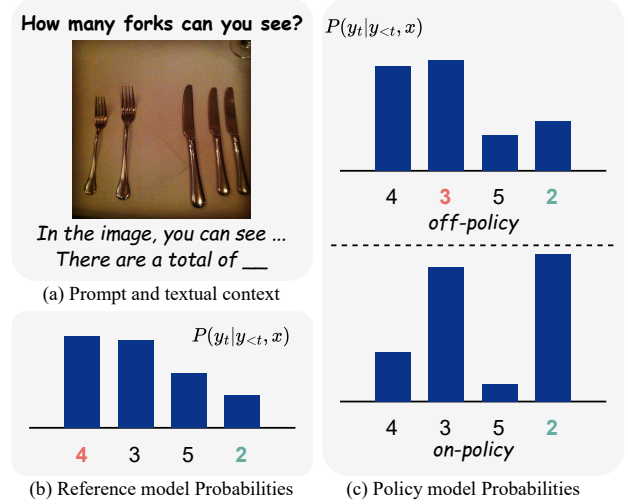(b) Reference model Probabilities

(c) Policy model Probabilities

Figure 1. An illustrative example of generation probability distribution after on/off-policy training. The correct token "2" is shown in green and the hallucinated token with the highest probability is shown in red. Panel (b) displays the hallucination mode of the reference model. Across the two training paradigms, off-policy training (panel (c), top) fails to overturn this dominant hallucination pattern. In contrast, on-policy training (panel (c), bottom) substantially increases the probability of the correct answer and effectively suppresses the dominant hallucinated token. The detailed analysis of this phenomenon is presented in Section 4.1.

which leverages preference data to align the model's behavior with human intentions. Current research is mainly centered on constructing high-quality multimodal preference data for this purpose. Some approaches generate multiple responses per prompt using existing LVLMs [1, 20, 24, 35], then employ expert models like GPTs [23, 24] to rank responses based on criteria like hallucination level [17]. Alternative approaches refine hallucinated components in model outputs to create positive samples [43, 47, 52], while others directly use ground-truth answers as positive samples and construct negatives by injecting controlled hallucinations [32]. Such data can be pre-processed and incorporated into offline datasets for subsequent model training. Conversely, some studies propose online data collection, where

preference data is generated on-the-fly during model update [35, 47, 49], which requires real-time data collection and annotation during training.

Among existing methods, most approaches construct data in an *off-policy* manner, where the data is generated from external models before training. In this work, we first analyze the training dynamics of preference alignment under the hallucination mitigation setting, during the analysis, we identify critical limitations in using off-policy data for LVLM hallucination mitigation, and we observe that on-policy data effectively addresses these limitations (Section 4.1), as shown in Figure 1. Therefore, we propose *an on-policy data construction pipeline*, and adopt an iterative updating mechanism, which is shown to surpass the performance of offline alternatives [36, 37, 45, 46].

In the on-policy paradigm, *reliable annotation* is key to data construction. To this end, there are two mainstream data annotation approaches: (1) Train a reward model for hallucination annotation [35]. While a local reward model reduces inference costs during annotation, training a reliably-evaluating reward model remains resource-intensive. (2) Leverage a fine-grained hallucination detector to rank samples [13, 49]. These approaches primarily assess the relative quality of model outputs based on the number of hallucinated segments or types. However, they focus exclusively on the relative quality, neglecting the potential intrinsic hallucinations that may still be present in supposedly superior responses. This method of modeling reward scores can, therefore, influence subsequent model optimization in ways that may not fully address underlying hallucinations. Overall, the inconvenience above leads to a question regarding reliable annotation:

*Are there better ways to provide annotation for on-policy data in hallucination mitigation?*

To address this question, we propose *hallucination-free chosen sample selection*, a novel data construction pipeline capable of generating high-quality on-policy data. In the algorithm level, we choose the online version of DPO– iterative DPO [45] as the preference alignment algorithm, and introduce a robust sample reweighting method by assigning higher weights to more informative pairs during training. In summary, our contributions are threefold:

1. We identify limitations of off-policy learning in addressing LVLM hallucination, and propose an effective pipeline to construct data in an on-policy manner.

2. We design an effective iterative alignment paradigm with a robust sample reweighting algorithm for training.

3. We conduct comprehensive experiments across multiple benchmarks and compare our approach with state-of-the-art baselines, empirically demonstrating the effectiveness and efficiency of our proposed method.

## 2. Related Works

**Hallucination Mitigation in LVLMs.** The hallucination phenomenon in LVLMs can originate from either the visual encoder [9, 11] or the pretrained LLM [15, 16]. These components may fail to fully align visual and textual representations, leading to inconsistencies in generated outputs. To address this issue, various visual encoders have been developed to enhance the quality of processed images, ensuring more accurate and contextually relevant outputs [1, 4, 50]. Additionally, fine-tuning LVLMs on datasets specifically curated to address hallucination has proven effective in enhancing alignment [35, 40, 53]. Another promising approach is contrastive decoding, which leverages the difference between image-conditioned and image-free token probabilities during decoding stage to prioritize tokens that are grounded in the visual information [6, 14, 16].

In this paper, we primarily focus on addressing hallucination of LVLMs through preference alignment, where it is critical to construct informative and high-quality preference pairs to guide the model in generating grounded responses. Numerous methods have been proposed for constructing offline hallucination preference datasets. These include contaminating or removing image content to create negative samples [27, 40, 44], injecting hallucination into textual responses to generate negative samples [32, 53], and leveraging human annotators or external expert models, such as GPTs, to refine generated responses and construct positive samples [43, 52]. Some works also propose constructing preference dataset in an on-policy manner [47, 49, 54].

**Preference Alignment.** Preference alignment has emerged as a cornerstone methodology for enhancing the response quality of LLMs [2, 5, 25, 39]. Central to this approach is reinforcement learning from human feedback (RLHF) [2, 25], which involves training a reward model to capture human preferences and then using reinforcement learning algorithms, such as Proximal Policy Optimization (PPO) [34], to guide LLMs toward generating responses with higher rewards. However, RL-based methods often face challenges related to instability during training. Consequently, recent research has shifted toward developing simpler and more stable alternatives to RLHF. A notable approach is DPO [28], which implicitly optimizes the same objective as RLHF but achieves human preference alignment through a single cross-entropy loss, bypassing the need for learning the explicit reward model and the complex reinforcement learning stage.

## 3. Preliminaries

**Large Vision Language Models.** Like LLMs, LVLMs operate by progressively predicting the probability distribution of the next token for a given prompt. In LVLMs, the
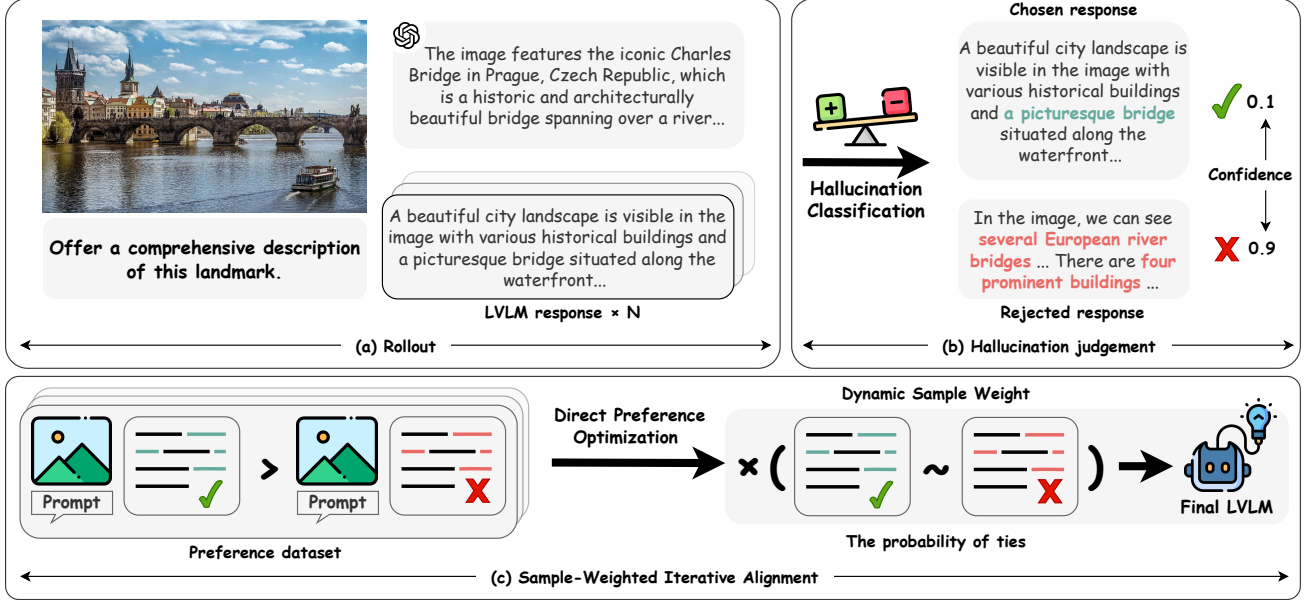
Figure 2. Overview of our framework. Our method consists of three steps: (1) Rollout: generating $N$ responses per image-prompt pair to form ⟨image, prompt, GT answer, response⟩ tuples; (2) Hallucination Judgement: selecting chosen and rejected samples based on hallucination probabilities from a trained classifier; (3) Sample-Weighted Iterative Alignment: fine-tuning the model using the preference dataset. These steps are repeated iteratively until the model converges.

prompt $x$ consists of a multimodal pair—an image $x_v$ and accompanying text $x_t$. For simplicity, we use $x$ to denote the unified prompt combining both visual and textual components ($x_v$ and $x_t$). Given this combined input, the LVLM then generates a textual response $y$, leveraging its ability to interpret and integrate information from both modalities.

Let an LVLM take an input $x \in \mathcal{X}$ and generate an output $y \in \mathcal{Y}$. In the contextual bandit formulation for RLHF, the LVLM is viewed as a policy $\pi_\theta(y \mid x)$ parameterized by $\theta$, which outputs an action $y$ (response) based on the state $x$ (prompt). Preference data is further collected and annotated by human labelers or AI feedback, denoted as $y_w \succ y_l \mid x$, where $y_w$ is the chosen response and $y_l$ is the rejected one in two responses generated for the prompt $x$.

**Direct Preference Optimization.** DPO is a direct preference alignment algorithm that unifies the reward learning and policy optimization stages in standard RLHF pipeline. Unlike traditional RLHF, DPO directly optimizes the policy on the offline dataset $D = \{x^{(i)}, y_w^{(i)}, y_l^{(i)}\}_{i=1}^N$ with $N$ preference samples, eliminating the need to learn an explicit reward model. Recall that in RLHF, we have the policy optimization objective as:

$$\mathcal{L}_\pi(\theta) = -\,\mathbb{E}_{x \sim P_x, y \sim \pi_\theta}\left[r_\phi(x, y)\right] + \beta \mathbb{D}_{KL}\left[\pi_\theta(y|x) || \pi_{\text{ref}}(y|x)\right], \quad (1)$$

where $P_x$ is the distribution of the prompt $x$, $r_\phi(\cdot)$ is the parameterized reward function, and $\pi_{\text{ref}}$ is the initial reference model.

As demonstrated by [26], the optimal conditional distribution $\pi_r(y \mid x)$ that minimizes the loss function in Equation (1) has the following closed-form solution:

$$\pi_r(y \mid x) = \frac{1}{Z(x)}\pi_{\text{ref}}(y \mid x)\exp\left(\frac{1}{\beta}r_\phi(x, y)\right), \quad (2)$$

where $Z(x) = \sum_y \pi_{\text{ref}}(y \mid x)\exp\left(1/\beta \cdot r_\phi(x, y)\right)$ is the partition function that guarantees normalization for $\pi_r$.

By leveraging this result, the DPO loss function is derived from the reward modeling objective, yielding a simplified optimization problem where only the parameterized policy $\pi_\theta$ serves as the optimization variable:

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(x, y_w, y_l) \sim D}\left[\log \sigma\left(\hat{r}(x, y_w) - \hat{r}(x, y_l)\right)\right], \quad (3)$$

where $\hat{r}(x, y) = \beta \log\left(\pi_\theta(y|x)/\pi_{\text{ref}}(y|x)\right)$ represents the implicit reward model derived from $\pi_\theta$.

## 4. Methodology

In this section, we first present two observations demonstrating the limitations of off-policy update for mitigating hallucinations in LVLMs. Then we propose an efficient data generation pipeline under the on-policy paradigm using a specifically-trained hallucination classifier. Finally, we present the detailed design of our iterative DPO algorithm with sample reweighting.

## 4.1. Key Observations

**Observation 1. DPO is a reweighting of the reference model.** According to the original objective of policy optimization in RLHF, the closed-form solution can be derived for the optimal policy $\pi_r(y \mid x)$, which is illustrated in Equation (2). We can generalize the expression following [7],

$$\pi_r(y \mid x) \propto f_x(r(x,y))\pi_{\text{ref}}(y \mid x), \qquad (4)$$

where $f_x(\cdot)$ is a non-decreasing function dependent on prompt $x$.

From Equation (4), the optimal policy $\pi_r(y \mid x)$ is exactly a reweighted version of the reference model $\pi_{\text{ref}}(y \mid x)$. As DPO utilizes this policy form for optimization, the training process would only modify the probabilities of responses that lie within the support of $\pi_{\text{ref}}$. Consider a response $y$ that is rarely generated by $\pi_{\text{ref}}$, i.e., $\pi_{\text{ref}}(y \mid x) \to 0$. Given that the weight $f_x(r(x,y))$ is bounded, even if $y \mid x$ appears in the training set as a chosen sample, the optimized generation probability for $y \mid x$ of the parameterized policy $\pi_\theta$ remains negligible, which means that $\pi_\theta(y \mid x) \to 0$.

This is often the case when using off-policy data for alignment, where the data distribution varies greatly from the model we are trying to optimize. This case becomes particularly relevant for the hallucination problem in LVLMs, where the positive responses are generally ground-truth answers, or GPTs modified responses, which are unlikely to be generated by the reference model. From the reweighting perspective, such phenomenon would substantially reduce the influence of chosen samples on the model's output distribution.

**Observation 2. Dominant hallucination patterns remain dominant even after off-policy training.** To rigorously characterize this observation, we begin by providing a formal definition of hallucination mitigation.

**Hallucination Mitigation.** Consider an LVLM with a vocabulary size of $\mathcal{V}$. When given an input token sequence prefix $\boldsymbol{x}$, the model predicts the distribution $\boldsymbol{p} \in \mathbb{R}^{\mathcal{V}}$ over the next token $z = x_{i+1}$ via its softmax output head. Let $\boldsymbol{p}^t$ represent the token distribution at training step $t$, with $\boldsymbol{p}_i^t$ denoting the corresponding probability for token $i$. Based on the above problem formulation, hallucination mitigation is formally characterized as a reallocation of probability mass from a hallucinatory token to a veridical alternative. Assume the model hallucinates, so that the most–likely token $h = \arg\max_{k \in \mathcal{V}} \boldsymbol{p}_k^t$ belongs to the hallucination set $\mathcal{H} \subset \mathcal{V}$. Define the set of non-hallucinatory tokens as $\mathcal{C} = \mathcal{V} \setminus \mathcal{H}$. The objective is to decrease the likelihood

of $h$ while increasing the likelihood of a corrective token $c \in \mathcal{C}$ until at a certain training step $t$, there exists

$$\boldsymbol{p}_c^t > \boldsymbol{p}_h^t. \qquad (5)$$

Once this inequality is satisfied, $c$ becomes the modal choice under greedy decoding, ensuring that the model emits a non-hallucinatory output.

The following remark analyzes the probability gap between the hallucinated choice $h$ and a non-hallucinated response $c$ after one gradient step on a single training pair $(\boldsymbol{x}, z)$.

*Remark* 4.1. Under the off-policy update paradigm, we have that for any $c \in \mathcal{C}$ and $c \neq z$,

$$\boldsymbol{p}_h^{t+1} - \boldsymbol{p}_c^{t+1} \geq \boldsymbol{p}_h^t - \boldsymbol{p}_c^t \geq 0. \qquad (6)$$

According to Equation (6), we can derive the inequality that $\boldsymbol{p}_h^{t+1} \geq \boldsymbol{p}_c^{t+1}$ (We refer the readers to Section B.2 for complete proof). By iteratively applying this inequality, we can have the result that after training, the relationship $\boldsymbol{p}_h \geq \boldsymbol{p}_c$ still holds for any $c \in \mathcal{C}$.

This phenomenon reveals a primary drawback of off-policy alignment in hallucination mitigation, that the hallucination pattern exhibited by the original model continues to take up high probability mass even after off-policy alignment. This observation is also evident in cases such as shown in Figure 1, where hallucinated tokens remains the highest probabilities after off-policy alignment, unlike the case with on-policy updates.

From both observations, we conclude that the on-policy paradigm demonstrates a distinct advantage over the off-policy approach in mitigating hallucinations. (This conclusion is also mentioned in Yang et al. [47], which gave a qualitative analysis of the gradient confined in DPO, whereas we provide a quantitative result concerning the output changes in all preference alignment algorithms). The superiority of the on-policy paradigm stems from its direct suppression of the model's inherent hallucination pattern and its capability to correct hallucination outputs into non-hallucinated ones, both of which are infeasible in off-policy alignment.

## 4.2. Hallucination-Free Chosen Sample Selection

Although on-policy learning is advantageous over its off-policy counterpart, a crucial problem emerges. In hallucination mitigation, it is vital to guarantee that *the chosen samples in training data do not contain any hallucinated contents*, due to the fact that the patterns in chosen samples are reinforced during preference alignment. This condition is naturally met in off-policy learning, where the chosen samples are generally constructed using ground-truth responses in pre-existing datasets [27, 53], or generated high-quality responses from expert models (like GPT-4) [43, 48].
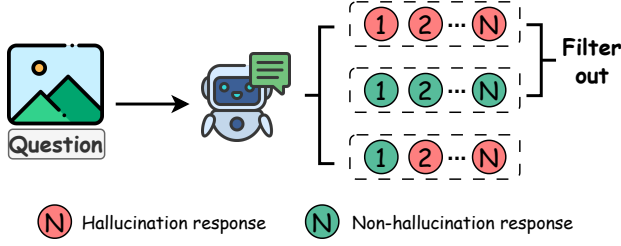
**Figure 3.** Hallucination-Free Chosen Sample Selection Process. For each prompt, we generate $N$ responses and use a hallucination classifier to evaluate each response. Prompts where all responses are either entirely hallucinated or entirely hallucination-free are excluded from the training set.

However, this requirement has not been handled properly in the on-policy paradigm, where both chosen and rejected data are sampled on-the-fly during training. The typical approach to construct training data in the on-policy setting involves adopting external models or rules to judge the samples' relative quality [13, 35, 49]. However, such annotation procedure does not guarantee that the chosen samples are free from hallucination.

To address this problem, we establish an evaluation criterion based on whether the sample contains hallucination or not. This criterion is then used to filter out hallucinated samples from the chosen. Specifically, a classifier is trained to label samples as either hallucinated or non-hallucinated for preference data annotation.

Note that our classifier differs fundamentally from existing classifier-based methods in that our approach evaluates hallucinations at the response level rather than performing sentence-by-sentence analysis, thereby dramatically reducing annotation costs. We further point out that the fine-grained classification approaches hinge on the assumption that the level of hallucination can be rated based on manually defined rules, such as the count of hallucinated elements within a sample [43], whereas we believe the semantics behind hallucination can hardly be captured by such simplistic predefined rules. Therefore we advocate for *learning hallucination patterns from data*, through training a comprehensive hallucination classifier.

To further improve the annotation accuracy of the hallucination classifier, we inject the ground-truth answer for each question as side information during the labeling process, which helps the model more accurately identify hallucinated content. As a result, our hallucination classifier takes the following input:

> **System prompt** $\backslash n$ [IMAGE_TOKEN] $\backslash n$
> **Question:** {}$\backslash n$ **Correct answer:** {}$\backslash n \backslash n$
> **Model response:** {}[EOS_TOKEN]

After the classifier is ready, we then construct training

data for DPO. Images, prompts, and ground truth answers are first collected from existing public datasets. As shown in panel (a) of Figure 2, for each prompt, we sample $N$ diverse responses from the model through appropriate randomization. Our trained classifier then evaluates each response's hallucination status, using threshold $\tau = 0.5$ to categorize samples into hallucinated/non-hallucinated categories, corresponding to panel (b) of Figure 2. We then select chosen and rejected samples for each prompt as:

- **Chosen:** Non-hallucinated sample with the lowest hallucination probability.
- **Rejected:** Hallucinated sample with the highest probability.

Note that we eliminate prompts where the sampled responses exclusively contain hallucinations or no hallucinations, thus avoiding situations where low-quality responses appear in chosen samples or high-quality responses appear in rejected samples, as shown in Figure 3.

### 4.3. Sample-Weighted Iterative Alignment

To further harness the power of on-policy data, we perform model update and data collection in an iterative manner, where multiple rounds of on-policy data are used. After each iteration, we gather new in-distribution responses from the updated model and apply our hallucination classifier to construct preference pairs for the subsequent training. During training, we observe that samples contribute unequally to the model's convergence. To optimize the learning process, we introduce a sample reweighting strategy. Our approach is informed by the implicit reward in DPO, which indicates the model's learning status for a given preference pair [42]. Specifically, we categorize samples based on their implicit reward margin:

- **Easy samples.** A large positive reward margin signifies that the model can already distinguish the preference pair with high confidence. These samples have been mastered and offer diminishing marginal returns during further training.
- **Hard samples.** A large negative reward margin indicates significant difficulty in distinguishing the preference. This may stem from a genuinely challenging case or potential annotation errors, where responses might be out-of-distribution for the classifier.
- **Boundary samples.** A reward margin close to zero implies high model uncertainty, implying substantial learning potential. We identify these as boundary samples, as the sample lies near the decision boundary of the current model.

Based on this categorization, we augment the DPO loss function to incorporate sample weights. We assign higher weights to boundary samples to maximize their learning

impact, while assigning lower weights to both easy and hard/noisy samples to improve robustness [19, 51].

We begin by introducing the Rao-Kupper model [29]. Unlike the BT model, which solely predicts the probability of wins/loses for a pair, the Rao-Kupper model also accounts for the probability of ties. Given a prompt $x$ and two responses $y_i$ and $y_j$, assuming the ground truth reward function is $r(\cdot)$, the model incorporates a parameter $\nu$ that governs the allocation of probability to ties.

Let $r_i = r(x, y_i)$ and $r_j = r(x, y_j)$, we have

$$
\begin{aligned}
p(y_i \succ y_j \mid x) &= \frac{1}{1 + \nu e^{(r_i - r_j)}}, \\
p(y_i \sim y_j \mid x) &= \frac{\nu^2 - 1}{\left(1 + \nu e^{(r_i - r_j)}\right)\left(1 + \nu e^{(r_j - r_i)}\right)},
\end{aligned}
\tag{7}
$$

where we set $\nu = 3.0$.

By incorporating this dynamic weight in the DPO loss, we obtain our loss for the iterative weighted DPO:

$$
\mathcal{L}_{\text{DPO}}^{\text{w}}(\theta) = -\mathbb{E}\left[\text{sg}\left(p(y_w \sim y_l \mid x) + \frac{2}{\nu + 1}\right) \cdot \ell_{\text{pair}}\right],
\tag{8}
$$

where $\ell_{\text{pair}} = \log \sigma\left(\hat{r}(x, y_w) - \hat{r}(x, y_l)\right)$, and $\text{sg}(\cdot)$ represents the stop-gradient operator. The expectation is taken over each batch of training data. Here, we add a bias term $2/(\nu + 1)$ to incorporate the original DPO loss ($\nu = 1$) into our framework.

The complete workflow of our robust iterative alignment algorithm is illustrated in Algorithm 1 in the Section E.

# 5. Experiments

In this section, we first empirically investigate the effectiveness of our method in addressing the hallucination problem for LVLMs. Then we analyze the efficacy of different components through ablation studies. We also give some cases of model generation to show the effect of our method.

## 5.1. Experimental Setup

**Models and Datasets.** For the classifier, we adapt Qwen2-VL-7B-Instruct [1] into a multimodal classifier by adding a linear projection head. To maintain consistency with previous work, we select LLaVA-1.5-7B and LLaVA-1.5-13B [20] as the base model for alignment. We obtain the training prompts from RLHF-V [48] by removing duplicate prompts, leaving 4.2k unique prompts. We also sample 3.4k distinct prompts from the VLFeedback dataset [17].

Our training procedure consists of one iteration of off-policy training followed by one iteration of on-policy training. In off-policy training, we designated the ground truth answer for each prompt as the chosen sample, while selecting model-generated response exhibiting hallucinations as the rejected sample. In the second iteration of on-policy

training, from the five candidate responses generated by our model, we construct the preference pair through data construction pipeline detailed in Section 4.2. At last, we utilize 6.4k prompts, each equipped with one pair of responses for the 7B and 13B models, respectively.

**Evaluation Metrics.** We conduct experiments on four multimodal benchmarks: (1) (1) AMBER [41] is a multi-dimensional benchmark for assessing LVLM hallucinations in generation and discrimination tasks. We use its generation subset of 1,004 questions and follow the standard protocol to report CHAIR, object coverage, hallucination rate, and human cognition alignment metrics. (2) MMHal-Bench [35] is a rigorously constructed benchmark for assessing multimodal hallucinations. It spans 12 COCO object categories and 8 task types, using GPT-4 scoring (0–6 scale) to compute average scores and hallucination rates. (3) Object HalBench [31] is a standard benchmark for evaluating hallucinations in captioning tasks. Models must describe 300 unique images in detail. Following RLAIF-V, we report sentence-level (CHAIRs, CHAIRsr) and object-level (CHAIRi) metrics. (4) MMBench [22] evaluates general multimodal understanding ability. We include this benchmark to complement our hallucination-focused evaluations; however, due to space limitations, the full results are deferred to the Table 5.

**Baseline Algorithms.** We compare our approach with state-of-the-art baselines. For base models, we select GPT-4V [24], Qwen-VL-Chat [1], LLaVA-1.5-7B, and LLaVA-1.5-13B [20]. The alignment methods used for comparison include LLaVA-RLHF [35], HALVA [32], mDPO [40], HA-DPO [52], POVID [53], RLAIF-V [49], OPA-DPO [47], RLHF-V [48], and HSA-DPO [43].

## 5.2. Main Results

We conduct quantitative experiments of our methods against baselines across various hallucination benchmarks. From Table 1, we draw the following conclusions:

- Our method achieves state-of-the-art performance on the vast majority of metrics across different hallucination benchmarks. For example, on MMHalBench, our approach significantly reduces the hallucination rates of LLaVA-1.5-7B and LLaVA-1.5-13B by 50.8% and 51.9%, respectively. For other metrics, such as CHAIRs and CHAIRi on Object HalBench, our 13B model achieves scores of 11.33 and 2.56, respectively, far surpassing existing methods. In summary, our approach demonstrats consistent performance across diverse hallucination benchmarks.

- Hallucination classifier exhibits high performance in data construction. Among the methods we compare, LLaVA-RLHF utilizes a high-quality human-annotated dataset containing 20k responses to train a 13B-sized reward model. However, the model trained with feed-

Table 1. Quantitative results of LLaVA-1.5-7B and LLaVA-1.5-13B trained with different preference optimization methods cross various hallucination benchmarks. For reference, we also provide additional results using various multimodal LLMs, preference data, and learning objectives, although these are not directly comparable. The best result for each metric within each group is highlighted in bold, and the second-best is underlined.

| Algorithm | Feedback | AMBER (1004) | | | | MMHal-Bench (96) | | Object Hal (300) | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | CHAIR ↓ | Cover ↑ | HalRate ↓ | Cog ↓ | Score ↑ | HalRate ↓ | CHAIRs ↓ | CHAIRsr ↓ | CHAIRi ↓ |
| **GPT-4V [24]** | | 4.6 | 67.1 | 30.7 | 2.6 | 3.49 | 0.28 | 13.6 | - | 7.3 |
| **Qwen-VL-Chat-34B [1]** | | 6.6 | 53.2 | 31.0 | 2.9 | 2.89 | 0.43 | 36 | - | 21.3 |
| +Silkie [17] *(EMNLP'24)* | GPT-4V | 5.4 | 55.8 | 29.0 | 2.0 | 3.01 | 0.41 | 25.3 | - | 13.9 |
| **LLaVA-1.5-7B [20, 21]** | | 7.7 | 51.6 | 34.7 | 4.2 | 2.01 | 0.61 | 55.00 | 55.18 | 16.02 |
| +LLaVA-RLHF [35] *(ACL'24)* | Reward-Model | 9.7 | **53.2** | 46.6 | 5.3 | 2.04 | 0.68 | 51.33 | 51.51 | 15.26 |
| +HALVA [32] *(ICLR'25)* | GPT-4V | 6.6 | <u>53.0</u> | 32.2 | 3.4 | 2.25 | 0.54 | 41.40 | - | 11.70 |
| +mDPO [40] *(EMNLP'24)* | GPT-4V | 4.4 | 52.4 | 24.5 | 2.4 | 2.39 | 0.54 | 35.70 | - | 9.80 |
| +HA-DPO [52] *(arXiv'23)* | GPT-4 | 7.8 | 52.1 | 35.6 | 4.2 | 1.89 | 0.65 | 53.33 | 53.33 | 9.58 |
| +POVID [53] *(arXiv'24)* | GPT-4V | 5.0 | 50.1 | 28.6 | 3.0 | 2.08 | 0.56 | 36.67 | 36.67 | 15.43 |
| +RLAIF-V [49] *(CVPR'25)* | LLaVA-Next | 3.0 | 50.4 | 16.2 | 1.0 | <u>3.00</u> | <u>0.38</u> | 14.67 | 14.81 | <u>3.83</u> |
| +OPA-DPO [47] *(CVPR'25)* | GPT-4V | <u>2.6</u> | 45.5 | **11.4** | <u>0.9</u> | 2.83 | 0.45 | <u>14.00</u> | 14.53 | 4.08 |
| **+Ours** | Qwen2-VL-7B | **2.4** | 48.6 | <u>13.6</u> | **0.8** | 3.01 | **0.30** | **12.33** | **12.67** | **2.99** |
| **LLaVA-1.5-13B [20, 21]** | | 6.8 | 51.9 | 31.8 | 3.3 | 2.48 | 0.52 | 52.00 | 52.17 | 14.46 |
| +LLaVA-RLHF [35] | Reward-Model | 7.7 | <u>52.3</u> | 38.6 | 4.0 | 2.53 | 0.57 | 47.33 | 47.65 | 13.21 |
| +RLHF-V (HD) [48] *(CVPR'24)* | Human | 6.3 | 46.1 | 25.1 | 2.1 | 2.81 | 0.49 | - | - | - |
| +HSA-DPO [43] | GPT-4/4V | <u>2.1</u> | 47.3 | 13.4 | 1.2 | 2.61 | 0.48 | - | - | - |
| +HALVA [32] | GPT-4V | 6.4 | **52.6** | 30.4 | 3.2 | 2.58 | 0.45 | 45.40 | - | 12.80 |
| +OPA-DPO [47] | GPT-4V | 2.6 | 48.0 | **12.6** | **1.0** | <u>3.07</u> | <u>0.39</u> | <u>16.00</u> | 16.16 | <u>4.87</u> |
| **+Ours** | Qwen2-VL-7B | **2.0** | 47.9 | <u>12.8</u> | **1.0** | **3.36** | **0.25** | **11.33** | **11.60** | **2.56** |

Table 2. Ablation study on sample reweighting (ObjHal.: Object HalBench; MMhHal.: MMHal-Bench).

| Data | ObjHal. | | MMhHal. | |
|---|---|---|---|---|
| | CHAIRs ↓ | CHAIRi. ↓ | Score ↑ | HalRate ↓ |
| iteration 1 | 48.67 | 14.55 | 2.53 | 0.50 |
| w/o reweight | 13.67 | 3.45 | 2.86 | 0.34 |
| Ours | **12.33** | **2.99** | **3.01** | **0.30** |

Table 3. Ablation study on on-policy data and filtering prompts (ObjHal.: Object HalBench; MMhHal.: MMHal-Bench).

| Data | ObjHal. | | MMhHal. | |
|---|---|---|---|---|
| | CHAIRs ↓ | CHAIRi. ↓ | Score ↑ | HalRate ↓ |
| w/o reweight | **13.67** | **3.45** | **2.86** | **0.34** |
| off-policy | 24.33 | 7.20 | 2.40 | 0.48 |
| w/o filtering | 19.00 | 4.82 | 2.14 | 0.65 |

back from this reward model did not effectively mitigate hallucinations. In contrast, we train a 7B hallucination classifier using only 8.4K responses annotated by DeepSeek-V3, achieving better results by a large margin. This result suggests that framing hallucination mitigation as a classification task yields superior results compared to reward model-based approaches.

• Our method is highly adaptable to models of different scales. Both 7B and 13B variants of our method achieve low hallucination rates and high performance scores across multiple benchmarks. Notably, our 13B model consistently outperforms its 7B counterpart in hallucination reduction across the three benchmarks - an improvement where some baseline methods actually show degraded performance with larger models.

## 5.3. Ablation Studies

To validate the efficacy of each component in our framework, we perform several ablation studies on two bench-marks. We first present the result of the first iteration in Table 2, to show that iterative online updates yield significant hallucination reduction, consistent with findings in [49]. We subsequently perform ablations on each of the three core components in our methodology: (1) sample reweighting strategy, (2) on-policy data training, and (3) reward model vs. classifier in data annotation.

**Sample reweighting.** For the sample reweighting ablation, we compare against a standard iterative DPO baseline with identical training data and hyperparameter configurations for fair comparison. As shown in Table 2, across various hallucination benchmarks, our method with the reweighting strategy consistently outperforms the standard DPO. This demonstrates that our approach can better balance the contributions of different samples during the training process and is more effectively aligned with the supervision signals provided by the classifier.

**On-policy data.** Further, we test the effectiveness of utilizing on-policy data against off-policy data. To eliminate the influence of other factors, we compare on-policy
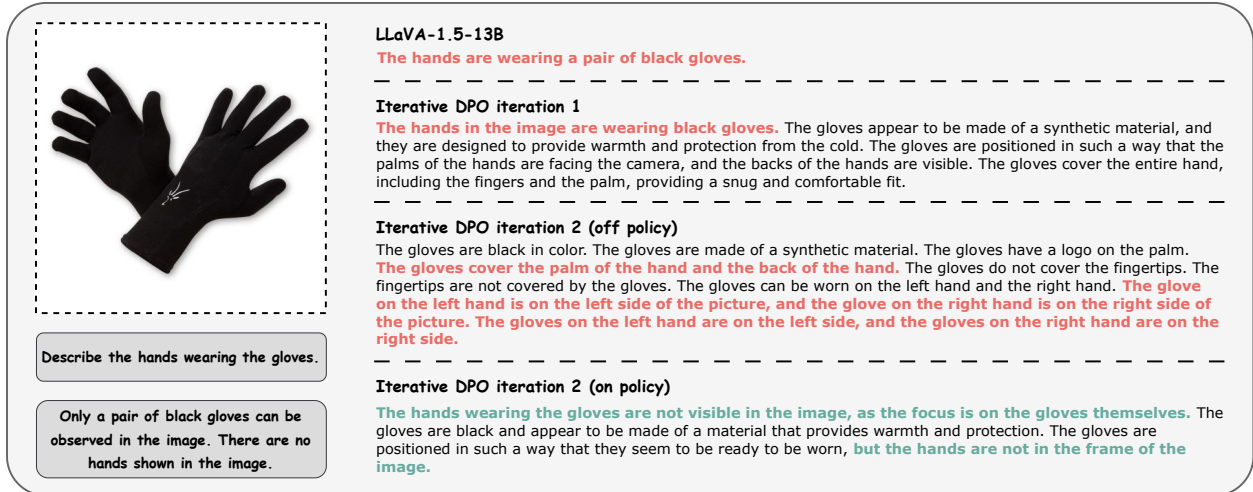
Figure 4. Qualitative Results from MMHalBench. We highlight the correct and incorrect parts of the responses from different models using bold green and red text, respectively.

and off-policy using the standard DPO loss with no sample reweighting. For the on-policy data utilized in our method, we sample and select chosen and rejected for each prompt as described in Section 4.2. To maintain comparable quality between two types of data, we replace chosen in the on-policy data with the ground truth answers from the corresponding dataset to construct off-policy data. The experimental results in Table 3 demonstrate that on-policy data can bring great improvements compared to off-policy data.

**Filtering hallucination-only and non-Hallucination Prompts.** To demonstrate that filtering out prompts where the sampled responses are either entirely hallucinated or completely non-hallucinated improves both training efficiency and reduces overall hallucination levels, we evaluate the model's performance on hallucination benchmarks. We compare the results with and without these filtered samples. Without filtering, the training data size nearly doubles. However, as shown in Table 3, better performance is achieved after filtering. This suggests that our method enhances the quality of the training data.

### 5.4. Case Study

To further illustrate the generation performance of our method, we present the responses from the base model LLaVA-1.5-13B, the two iterations of our method, and the alternative of iteration 2 of off-policy, respectively, to show the improvement on generation quality of our method in Figure 4. Except for on-policy iterative DPO iteration 2, which successfully answers the question correctly, the other results fail to explicitly identify the misleading nature inherent in the question itself. Off-policy iterative DPO iteration 1, while unable to resolve the model's hallucination problem, significantly increases response length – the key

motivation for this training phase. If we continue to adopt off-policy training in the iteration, it fails to address hallucinations and introduces repetitive content generation as shown in Figure 4. In contrast, our method not only highlights the misleading information within the question but also provides a detailed description of the image content.

## 6. Conclusion

In this work, we address the critical challenge of hallucination mitigation in LVLMs. We first recognize the superiority of on-policy data over off-policy alternatives for DPO in broad hallucination mitigation tasks, which highlights the necessity of collecting on-policy data for both positive and negative samples. This observation motivates us to propose an effective and efficient preference data construction pipeline to provide high-quality feedback for generated on-policy responses, via *training a binary hallucination classifier*. Experiments show that our method exhibits superior performance compared to mainstream baselines, demonstrating the advantage of adopting a classifier as critic for data annotation. To further enhance the power of on-policy data, we propose a robust iterative DPO algorithm to iteratively collect data and update policy. This approach employs a probabilistic model on DPO's implicit reward to dynamically determine sample weights, enabling adaptive prioritization of different sample types based on the model's learning condition. Our extensive evaluations across multiple benchmarks show consistent superiority in generating high-fidelity, low-hallucination responses across diverse image inputs. The method demonstrates strong performance on both LLaVA-1.5-7B and LLaVA-1.5-13B architectures, confirming its general applicability.

# References

[1] J Bai, S Bai, S Yang, S Wang, S Tan, P Wang, J Lin, C Zhou, and J Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. arxiv 2023. *arXiv preprint arXiv:2308.12966*, 2023. 1, 2, 6, 7, 12

[2] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022. 2, 12

[3] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018. 13

[4] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24185–24198, 2024. 2, 12

[5] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 1, 2, 12

[6] Alessandro Favero, Luca Zancato, Matthew Trager, Siddharth Choudhary, Pramuditha Perera, Alessandro Achille, Ashwin Swaminathan, and Stefano Soatto. Multi-modal hallucination control by visual information grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14303–14312, 2024. 2, 12

[7] Lin Gui, Cristina Gârbacea, and Victor Veitch. Bonbon alignment for large language models and the sweetness of best-of-n sampling. *arXiv preprint arXiv:2406.00832*, 2024. 4

[8] Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, et al. Direct language model alignment from online ai feedback. *arXiv preprint arXiv:2402.04792*, 2024. 12

[9] Xin He, Longhui Wei, Lingxi Xie, and Qi Tian. Incorporating visual experts to resolve the information loss in multimodal large language models. *arXiv preprint arXiv:2401.03105*, 2024. 2, 12

[10] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 2023. 1

[11] Jitesh Jain, Jianwei Yang, and Humphrey Shi. Vcoder: Versatile vision encoders for multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27992–28002, 2024. 2, 12

[12] Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on learning theory*, pages 2137–2143. PMLR, 2020. 13

[13] Liqiang Jing and Xinya Du. FGAIF: Aligning large vision-language models with fine-grained AI feedback. *Transactions on Machine Learning Research*, 2025. 2, 5

[14] Sihyeon Kim, Boryeong Cho, Sangmin Bae, Sumyeong Ahn, and Se-Young Yun. VACode: Visual augmented contrastive decoding. In *Trustworthy Multi-modal Foundation Models and AI Agents (TiFA)*, 2024. 2, 12

[15] Seongyun Lee, Sue Hyun Park, Yongrae Jo, and Minjoon Seo. Volcano: Mitigating multimodal hallucination through self-feedback guided revision. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 391–404, Mexico City, Mexico, 2024. Association for Computational Linguistics. 2, 12

[16] Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13872–13882, 2024. 2, 12

[17] Lei Li, Zhihui Xie, Mukai Li, Shunian Chen, Peiyi Wang, Liang Chen, Yazheng Yang, Benyou Wang, Lingpeng Kong, and Qi Liu. 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, 2024. Association for Computational Linguistics. 1, 6, 7

[18] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13*, pages 740–755. Springer, 2014. 18

[19] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 6

[20] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023. 1, 6, 7, 20

[21] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024. 1, 7, 20

[22] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, pages 216–233. Springer, 2024. 6

[23] OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. 1

[24] OpenAI. Gpt-4v(ision) system card. 2023. 1, 6, 7

[25] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini

Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022. 2, 12

[26] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019. 3

[27] Renjie Pi, Tianyang Han, Wei Xiong, Jipeng Zhang, Runtao Liu, Rui Pan, and Tong Zhang. Strengthening multimodal large language model with bootstrapped preference optimization. In *European Conference on Computer Vision*, pages 382–398. Springer, 2024. 2, 4, 12

[28] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023. 2, 12

[29] Pejaver V Rao and Lawrence L Kupper. Ties in paired-comparison experiments: A generalization of the bradley-terry model. *Journal of the American Statistical Association*, 62(317):194–204, 1967. 6

[30] Yi Ren and Danica J. Sutherland. Learning dynamics of LLM finetuning. In *The Thirteenth International Conference on Learning Representations*, 2025. 13

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

[32] Pritam Sarkar, Sayna Ebrahimi, Ali Etemad, Ahmad Beirami, Sercan O Arik, and Tomas Pfister. Mitigating object hallucination in MLLMs via data-augmented phrase-level alignment. In *The Thirteenth International Conference on Learning Representations*, 2025. 1, 2, 6, 7, 12

[33] Antoine Scheid, Etienne Boursier, Alain Durmus, Michael I Jordan, Pierre Ménard, Eric Moulines, and Michal Valko. Optimal design for reward modeling in rlhf. *arXiv preprint arXiv:2410.17055*, 2024. 13

[34] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 2, 12

[35] 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. Aligning large multimodal models with factually augmented RLHF. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 13088–13110, Bangkok, Thailand, 2024. Association for Computational Linguistics. 1, 2, 5, 6, 7, 12, 20

[36] Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Stefano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage suboptimal, on-policy data. In *ICML*, 2024. 2

[37] Yunhao Tang, Daniel Zhaohan Guo, Zeyu Zheng, Daniele Calandriello, Yuan Cao, Eugene Tarassov, Rémi Munos,

Bernardo Ávila Pires, Michal Valko, Yong Cheng, et al. Understanding the performance gap between online and offline alignment algorithms. *arXiv preprint arXiv:2405.08448*, 2024. 2

[38] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. 1

[39] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 2, 12

[40] Fei Wang, Wenxuan Zhou, James Y. Huang, Nan Xu, Sheng Zhang, Hoifung Poon, and Muhao Chen. mDPO: Conditional preference optimization for multimodal large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8078–8088, Miami, Florida, USA, 2024. Association for Computational Linguistics. 2, 6, 7, 12

[41] Junyang Wang, Yuhang Wang, Guohai Xu, Jing Zhang, Yukai Gu, Haitao Jia, Jiaqi Wang, Haiyang Xu, Ming Yan, Ji Zhang, et al. Amber: An llm-free multi-dimensional benchmark for mllms hallucination evaluation. *arXiv preprint arXiv:2311.07397*, 2023. 6

[42] Junkang Wu, Yuexiang Xie, Zhengyi Yang, Jiancan Wu, Jinyang Gao, Bolin Ding, Xiang Wang, and Xiangnan He. $\beta$-dpo: Direct preference optimization with dynamic $\beta$. *Advances in Neural Information Processing Systems*, 37: 129944–129966, 2024. 5

[43] Wenyi Xiao, Ziwei Huang, Leilei Gan, Wanggui He, Haoyuan Li, Zhelun Yu, Fangxun Shu, Hao Jiang, and Linchao Zhu. Detecting and mitigating hallucination in large vision language models via fine-grained ai feedback. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 25543–25551, 2025. 1, 2, 4, 5, 6, 7, 12

[44] Yuxi Xie, Guanzhen Li, Xiao Xu, and Min-Yen Kan. V-DPO: Mitigating hallucination in large vision language models via vision-guided direct preference optimization. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 13258–13273, Miami, Florida, USA, 2024. Association for Computational Linguistics. 2, 12

[45] Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative preference learning from human feedback: Bridging theory and practice for RLHF under KL-constraint. In *Proceedings of the 41st International Conference on Machine Learning*, pages 54715–54754. PMLR, 2024. 2

[46] Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenjie Ye, Weilin Liu, Zhiyu Mei, Guangju Wang, Chao Yu, and Yi Wu. Is dpo superior to ppo for llm alignment? a comprehensive study. In *ICML*, 2024. 2

[47] Zhihe Yang, Xufang Luo, Dongqi Han, Yunjian Xu, and Dongsheng Li. Mitigating hallucinations in large vision-language models via dpo: On-policy data hold the key. In *Proceedings of the Computer Vision and Pattern Recogni-*

*tion Conference*, pages 10610–10620, 2025. 1, 2, 4, 6, 7, 12, 20

[48] Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun, et al. Rlhf-v: Towards trustworthy mllms via behavior alignment from fine-grained correctional human feedback. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13807–13816, 2024. 1, 4, 6, 7, 12

[49] Tianyu Yu, Haoye Zhang, Qiming Li, Qixin Xu, Yuan Yao, Da Chen, Xiaoman Lu, Ganqu Cui, Yunkai Dang, Taiwen He, et al. Rlaif-v: Open-source ai feedback leads to super gpt-4v trustworthiness. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19985–19995, 2025. 1, 2, 5, 6, 7, 12, 20

[50] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023. 2, 12

[51] Zhilu Zhang and Mert R. Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, page 8792–8802, Red Hook, NY, USA, 2018. Curran Associates Inc. 6

[52] Zhiyuan Zhao, Bin Wang, Linke Ouyang, Xiaoyi Dong, Jiaqi Wang, and Conghui He. Beyond hallucinations: Enhancing lvlms through hallucination-aware direct preference optimization. *arXiv preprint arXiv:2311.16839*, 2023. 1, 2, 6, 7, 12, 20

[53] Yiyang Zhou, Chenhang Cui, Rafael Rafailov, Chelsea Finn, and Huaxiu Yao. Aligning modalities in vision large language models via preference fine-tuning. *arXiv preprint arXiv:2402.11411*, 2024. 1, 2, 4, 6, 7, 12, 15, 20

[54] Yiyang Zhou, Zhiyuan Fan, Dongjie Cheng, Sihan Yang, Zhaorun Chen, Chenhang Cui, Xiyao Wang, Yun Li, Linjun Zhang, and Huaxiu Yao. Calibrated self-rewarding vision language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 2, 12

# Supplementary Materials

## A. Related Works

**Hallucination Mitigation in LVLMs**  The hallucination phenomenon in LVLMs can originate from either the visual encoder [9, 11] or the pretrained LLM [15, 16]. These components may fail to fully align visual and textual representations, leading to inconsistencies in generated outputs. To address this issue, various visual encoders have been developed to enhance the quality of processed images, ensuring more accurate and contextually relevant outputs [1, 4, 50]. Additionally, fine-tuning LVLMs on datasets specifically curated to address hallucination has proven effective in enhancing alignment [35, 40, 53]. Another promising approach is contrastive decoding, which leverages the difference between image-conditioned and image-free token probabilities during decoding stage to prioritize tokens that are grounded in the visual information [6, 14, 16].

In this paper, we primarily focus on addressing hallucination of LVLMs through preference alignment, where it is critical to construct informative and high-quality preference pairs to guide the model in generating grounded responses. Numerous methods have been proposed for constructing offline hallucination preference datasets. These include contaminating or removing image content to create negative samples [27, 40, 44], injecting hallucination into textual responses to generate negative samples [32, 53], and leveraging human annotators or external expert models, such as GPTs, to refine generated responses and construct positive samples [43, 52]. Some works also propose constructing preference dataset in an on-policy manner [47, 49, 54].

**Preference Alignment.**  Preference alignment has emerged as a cornerstone methodology for enhancing the response quality of LLMs [2, 5, 25, 39]. Central to this approach is reinforcement learning from human feedback (RLHF) [2, 25], which involves training a reward model to capture human preferences and then using reinforcement learning algorithms, such as Proximal Policy Optimization (PPO) [34], to guide LLMs toward generating responses with higher rewards. However, RL-based methods often face challenges related to instability during training. Consequently, recent research has shifted toward developing simpler and more stable alternatives to RLHF. A notable approach is DPO [28], which implicitly optimizes the same objective as RLHF but achieves human preference alignment through a single cross-entropy loss, bypassing the need for learning the explicit reward model and the complex reinforcement learning stage.

The simplicity of DPO has inspired a wave of subsequent alternatives for hallucination mitigation in LVLMs [40, 48, 49, 52, 53]. Corresponding to how the dataset is constructed, the DPO algorithm can be tailored to address specific alignment challenges. For instance, some approaches focus on fine-grained preference feedback, enabling more nuanced alignment by capturing segment-level hallucination in responses [32, 43, 48]. Other than alignment on offline dataset, on-policy DPO [49] or its alternatives [47] emphasize aligning the model on its own generated outputs rather than the offline dataset. Furthermore, iterative DPO [49, 54] introduces an iterative updating paradigm similar to the standard RL process, progressively improving alignment over multiple iterations.

## B. Theoretical Proof

### B.1. Definition of on/off-policy in Preference Alignment

As stated in [8], the key distinction between on-policy and off-policy lies in whether the training data used to optimize the current policy is generated by the current policy itself. If the data is generated by the current policy, it is considered on-policy; otherwise, it is off-policy.

For a given completion $y$, if it is collected in an on-policy manner, i.e., $y \sim \pi_\theta(\cdot|x)$, then the current policy has a higher probability of generating that completion. In contrast, if the distribution that generated $y$ differs from the current policy (which, in general, is a significant difference), the probability of the current policy generating $y$ is relatively low, potentially approaching zero.

### B.2. Proof of Remark 4.1

*Proof.* The next-token prediction task is typically trained via maximum likelihood estimation (MLE), which is equivalent to minimizing the cross-entropy loss. Let $\boldsymbol{x} = (x_1, x_2, \ldots, x_i)$ denote a token sequence prefix, and let $y = x_{i+1}$ denote the next token to be predicted. A large language model with vocabulary size $\mathcal{V}$ maps $\boldsymbol{x}$ to a $d$-dimensional feature vector $\phi(\boldsymbol{x}) \in \mathbb{R}^d$ via a deep neural network. The model then computes a logit vector $\boldsymbol{z} = \boldsymbol{W}^\top \phi(\boldsymbol{x}) \in \mathbb{R}^{\mathcal{V}}$ using a linear transformation $\mathbf{w} \in \mathbb{R}^{d \times \mathcal{V}}$, and applies the Softmax function to obtain the predicted probability vector $\boldsymbol{p} = \mathrm{Softmax}(\boldsymbol{z})$. The standard cross-entropy loss is used during training, and it is defined as:

$$\mathcal{L}_{CE}(\boldsymbol{p}, y) = -\boldsymbol{e}_y^\top \log \boldsymbol{p}, \tag{9}$$

where $\boldsymbol{e}_y \in \mathbb{R}^{\mathcal{V}}$ is a one-hot vector with a $1$ at the $y$-th position and zeros elsewhere. We operate under the assumption of a linearly parametrized softmax policy [12, 30, 33], in which the feature extractor $\phi$ is fixed, and only the parameters of the read-out layer $\boldsymbol{W}$ are updated. Using stochastic gradient descent with learning rate $\eta$, the update rule is:

$$\boldsymbol{W}^{t+1} = \boldsymbol{W}^t - \eta \nabla_{\boldsymbol{W}} \mathcal{L}_{CE} = \boldsymbol{W}^t - \eta \phi(\boldsymbol{x})(\boldsymbol{p}^t - \boldsymbol{e}_y)^\top, \tag{10}$$

where $t$ denotes the $t$-th training step.

To analyze the dynamics of the training process, we convert the discrete update into a continuous-time differential equation [3]. Let $\boldsymbol{W}(t)$ denote the parameters at continuous time $t$, then:

$$\frac{d\boldsymbol{W}(t)}{dt} = -\eta \phi(\boldsymbol{x}) \left(\boldsymbol{p}(t) - \boldsymbol{e}_y\right)^\top. \tag{11}$$

Let $\boldsymbol{z}(t) = \boldsymbol{W}(t)^\top \phi(\boldsymbol{x})$, then $\boldsymbol{p}(t) = \mathrm{Softmax}(\boldsymbol{z}(t))$. Differentiating $\boldsymbol{p}(t)$ with respect to time yields:

$$\frac{d\boldsymbol{p}(t)}{dt} = \frac{d \, \mathrm{Softmax}(\boldsymbol{z}(t))}{d\boldsymbol{z}(t)} \cdot \frac{d\boldsymbol{z}(t)}{dt}. \tag{12}$$

Since $\phi(\boldsymbol{x})$ is constant, we have:

$$\frac{d\boldsymbol{z}(t)}{dt} = \left(\frac{d\boldsymbol{W}(t)}{dt}\right)^\top \phi(\boldsymbol{x}) = -\eta \left[\phi(\boldsymbol{x})^\top \phi(\boldsymbol{x})\right] (\boldsymbol{p}(t) - \boldsymbol{e}_y). \tag{13}$$

Let $\beta = \eta \|\phi(\boldsymbol{x})\|^2$ for convenience. On the other hand, the Jacobian matrix of the Softmax function is:

$$\frac{d\boldsymbol{p}}{d\mathbf{z}} = \mathrm{diag}(\boldsymbol{p}) - \boldsymbol{p}\boldsymbol{p}^\top. \tag{14}$$

Substituting these results gives the time derivative of $\boldsymbol{p}(t)$:

$$\frac{d\boldsymbol{p}(t)}{dt} = -\beta \left[\mathrm{diag}(\boldsymbol{p}(t)) - \boldsymbol{p}(t)\boldsymbol{p}(t)^\top\right] (\boldsymbol{p}(t) - \mathbf{e}_y). \tag{15}$$

This is a nonlinear vector differential equation describing the evolution of the prediction probabilities $\boldsymbol{p}(t)$ under gradient descent training. To extract the dynamics of each component $\boldsymbol{p}_k(t)$, we expand the above expression as:

$$\frac{d\boldsymbol{p}_k(t)}{dt} = -\beta \boldsymbol{p}_k(t) \left[ (\boldsymbol{p}_k(t) - \delta_{ky}) - \sum_{j=1}^{V} \boldsymbol{p}_j(t)(\boldsymbol{p}_j(t) - \delta_{jy}) \right], \tag{16}$$

where $\delta_{ky}$ is the Kronecker delta. When $k = y$, $\delta_{ky} = 1$; otherwise, it is $0$.

To numerically solve this continuous-time system, we apply the Euler method. Let the time step be $\Delta t$, and define $\boldsymbol{p}^{(n)} := \boldsymbol{p}(t_n)$ at discrete time $t_n = n\Delta t$. The Euler update rule is:

$$\boldsymbol{p}^{(n+1)} = \boldsymbol{p}^{(n)} - \Delta t \cdot \beta \left[ \operatorname{diag}(\boldsymbol{p}^{(n)}) - \boldsymbol{p}^{(n)} \boldsymbol{p}^{(n)\top} \right] (\boldsymbol{p}^{(n)} - \boldsymbol{e}_y), \tag{17}$$

and for the $k$-th component:

$$\boldsymbol{p}_k^{(n+1)} = \boldsymbol{p}_k^{(n)} - \Delta t \cdot \beta \cdot \boldsymbol{p}_k^{(n)} \left[ (\boldsymbol{p}_k^{(n)} - \delta_{ky}) - \sum_{j=1}^{\mathbf{V}} \boldsymbol{p}_j^{(n)}(\boldsymbol{p}_j^{(n)} - \delta_{jy}) \right]. \tag{18}$$

We analyze the relative dynamics of the hallucination component during training. Let $\boldsymbol{p}^t \in \mathbb{R}^{\mathcal{V}}$ denote the predicted probability vector at training step $t$. Let $y$ denote the ground-truth label and define the hallucination component as the non-ground-truth class with the highest probability:

$$\boldsymbol{p}_h^t := \max_{k \neq y} \boldsymbol{p}_k^t. \tag{19}$$

Let $c$ denote an arbitrary non-hallucination and non-target class, i.e., $c \notin \{h, y\}$. We show that $\boldsymbol{p}_h^{t+1} > \boldsymbol{p}_c^{t+1}$ always holds during training under the continuous-time Euler approximation of gradient descent. As derived from Equation (18), and omitting the superscript $(n)$ for clarity, we have:

$$\boldsymbol{p}_h^{(n+1)} = \boldsymbol{p}_h - \Delta t \cdot \beta \cdot \boldsymbol{p}_h \left[ \boldsymbol{p}_h - \sum_{j=1}^{\mathcal{V}} \boldsymbol{p}_j(\boldsymbol{p}_j - \delta_{jy}) \right]. \tag{20}$$

Define the auxiliary quantity:

$$f := \sum_{j=1}^{\mathcal{V}} \boldsymbol{p}_j(\boldsymbol{p}_j - \delta_{jy}) = \|\boldsymbol{p}\|^2 - \boldsymbol{p}_y. \tag{21}$$

Let $s := \boldsymbol{p}_h + \boldsymbol{p}_y$ be the total probability mass on the hallucination and target components, and define the residual mass:

$$R := 1 - s = \sum_{k \notin \{h,y\}} \boldsymbol{p}_k. \tag{22}$$

By the definition of $\boldsymbol{p}_h = \max_{k \neq y} \boldsymbol{p}_k$, we have for any $k \notin \{h, y\}$:

$$\sum_{k \notin \{h,y\}} \boldsymbol{p}_k^2 \leq \boldsymbol{p}_h \sum_{k \notin \{h,y\}} \boldsymbol{p}_k = \boldsymbol{p}_h R. \tag{23}$$

Thus, we can bound $f$ from below as follows:

$$
\begin{aligned}
f &= \boldsymbol{p}_h + \boldsymbol{p}_y - \left( \boldsymbol{p}_h^2 + \boldsymbol{p}_y^2 + \sum_{k \notin \{h,y\}} \boldsymbol{p}_k^2 \right) \\
&\geq \boldsymbol{p}_h + \boldsymbol{p}_y - \left( \boldsymbol{p}_h^2 + \boldsymbol{p}_y^2 + \boldsymbol{p}_h R \right) \\
&= \boldsymbol{p}_h + \boldsymbol{p}_y - \boldsymbol{p}_h^2 - \boldsymbol{p}_y^2 - \boldsymbol{p}_h(1 - \boldsymbol{p}_h - \boldsymbol{p}_y) \\
&= \boldsymbol{p}_y(1 - \boldsymbol{p}_y + \boldsymbol{p}_h) \geq 0.
\end{aligned}
\tag{24}
$$

Therefore, the hallucination probability satisfies

$$
\boldsymbol{p}_h \geq \|\boldsymbol{p}\|^2 - \boldsymbol{p}_y.
\tag{25}
$$

Next, we consider the update of a non-hallucination component $p_c$ for $c \notin \{h, y\}$. Its update is given by:

$$
\boldsymbol{p}_c^{(n+1)} = \boldsymbol{p}_c - \Delta t \cdot \beta \cdot \boldsymbol{p}_c \left[ \boldsymbol{p}_c - \left( \|\boldsymbol{p}\|^2 - \boldsymbol{p}_y \right) \right].
\tag{26}
$$

We now examine the difference between the updated hallucination and non-hallucination components:

$$
\boldsymbol{p}_h^{(n+1)} - \boldsymbol{p}_c^{(n+1)} = (\boldsymbol{p}_h - \boldsymbol{p}_c) - \beta \left[ \boldsymbol{p}_h(\boldsymbol{p}_h - f) - \boldsymbol{p}_c(\boldsymbol{p}_c - f) \right],
\tag{27}
$$

where $f = \|\boldsymbol{p}\|^2 - \boldsymbol{p}_y$.

Define the auxiliary function $g(x) := x(x - f)$, a quadratic function in $x$. Note that since $\boldsymbol{p}_h \geq f$, we have $g(\boldsymbol{p}_h) \geq 0$. Furthermore:
- If $\boldsymbol{p}_c \geq f$, then $g$ is increasing on $[f, 1]$, and $\boldsymbol{p}_h \geq \boldsymbol{p}_c$ implies $g(\boldsymbol{p}_h) \geq g(\boldsymbol{p}_c)$;
- If $\boldsymbol{p}_c < f$, then $\boldsymbol{g}(p_c) < 0 \leq g(\boldsymbol{p}_h)$.

In both cases, we conclude that

$$
g(\boldsymbol{p}_h) \geq g(\boldsymbol{p}_c),
\tag{28}
$$

which implies

$$
\boldsymbol{p}_h^{(n+1)} - \boldsymbol{p}_c^{(n+1)} \geq \boldsymbol{p}_h^{(n)} - \boldsymbol{p}_c^{(n)} \geq 0.
\tag{29}
$$

This indicates that function $d(t) = \boldsymbol{p}_h^t - \boldsymbol{p}_c^t$ is non-decreasing at any training step $t$. Moreover, since $d(t) \geq 0$, it follows that $d(t + 1) \geq 0$.

$\square$

## C. Experiment Setup

In this section, we present the complete experimental configuration, including implementation details and parameter specifications.

### C.1. Preparation of Classifier Training Data

We first construct a labeled dataset of model-generated samples with hallucination annotations for training the classifier. By incorporating ground truth annotations as auxiliary information, we simplify the classification task, enabling the model to make accurate judgments even when relying solely on the textual modality.

For classifier training, we extract data from the POVID dataset [53], using prompts as questions and chosen responses as ground truth answers, while using LLaVA-1.5-7B's outputs as model responses. For cost efficiency, we employ the pure text model DeepSeek-V3 for hallucination annotation. The detailed system prompt for annotation are listed as follows.

Please act as an impartial and objective judge and evaluate the quality of the response provided by a
Large Multimodal Model (LMM) to the user question. Your evaluation should be mainly based on whether
the response is informative, and whether the response contains any hallucination. Hallucination, in
this context, refers to a situation where the LMM generates a response that includes information not
present or implied in the image or previous conversation. A hallucination could be a false claim about
an object, action, emotion, or any other detail that is not grounded in the image.

For clarity, consider these examples:

### Example 1:
- Question: How many cats can you see in this image?
- Standard Human-Generated Answer: There are two cats visible in the image.
- LMM Response to Evaluate: There is one cat in the image.
- Analysis: The LMM's response is false because it claims there is only one cat in the image, when in
fact there are two.
- Hallucination: Yes.

### Example 2:
- Question: What can you see from this photo?
- Standard Human-Generated Answer: In this photo we can see branches. We can see a panda on a branch.
The background is blurry.
- LMM Response to Evaluate: In the photo, I can see a large black and white panda sitting on a tree
branch. The image captures the panda's unique appearance and its natural habitat, showcasing the beauty
 and rarity of these animals.
- Analysis: The LMM's response is not false because it does not make any false claims about the image
contents. For example, it does not claim that the panda is eating bamboo.
- Hallucination: No.

### Example 3:
- Question: Describe this image in details.
- Standard Human-Generated Answer: In this picture I can see a dog laying under a yellow umbrella on a
beach, and there is a rope on the right side. In the background I can see few persons talking or lying.
 At the top of the image, I can see the blue and clear sky.
- LMM Response to Evaluate: The image features a brown dog laying on the sand at the beach, underneath
a yellow umbrella. The dog is positioned near the water, and there are people in the background. The
beach is sandy, and the sun is shining, creating a pleasant atmosphere. The dog is enjoying the shade
provided by the umbrella, which is a common way to protect pets from the sun's harmful rays. The
presence of people in the background suggests that the beach is a popular spot for leisure activities,
such as sunbathing, swimming, and relaxing.
- Analysis: The LMM's response is false because it claims that the dog is positioned near the water,
but there is no water in the image.
- Hallucination: Yes.

### Example 4:
- Question: Where is this photo taken?
- Standard Human-Generated Answer: From the slogan "Welcome to Houston, Texas", we can infer that this
photo is taken at an airport in Houston, Texas.
- LMM Response to Evaluate: This photo is taken at an airport, specifically at the entrance of a
terminal. However, the specific location of this airport cannot be determined from the image.
- Analysis: Although the LMM's response does not explicitly mention the location "Houston, Texas" and
does not perfectly answer the question, it makes no false claims about the image contents. Therefore,
this response is not considered hallucination. It is OK if the LMM fails to directly answer the
question, expresses uncertainty, or cannot include all the objects or details in the human-generated
answer.
- Hallucination: No.

### Example 5:
- Question: What environmental concerns are involved in this photo?
- Standard Human-Generated Answer: This photo shows a smokestack emitting smoke. It may cause severe
air pollution which threatens human health.
- LMM Response to Evaluate: In the photo, there are several environmental concerns related to the
smokestack emitting smoke. The smoke from the smokestack is a byproduct of industrial processes, which
can contribute to air pollution, climate change, and negative impacts on human health. The smoke
contains harmful pollutants, such as particulate matter, sulfur dioxide, nitrogen oxides, and carbon
monoxide, which can lead to respiratory problems, heart disease, and other health issues. Additionally,
 the smoke contributes to the greenhouse effect, which can lead to global warming and climate change.

```
The smokestack's emissions also affect the environment, as they can harm wildlife, vegetation, and
ecosystems. Therefore, it is essential to address these environmental concerns by implementing measures
 to reduce emissions and promote sustainable practices in industrial processes.
- Analysis: Although the LMM's response is significantly longer than the standard human-generated
answer, it does not contain any false claims about the image contents. Instead, it provides additional
general information about the environmental concerns, which can be inferred from the smoke emission.
Such detailed analysis or reasoning should be considered as a positive aspect, as long as it contains
no false claims.
- Hallucination: No.

With these examples in mind, please help me evaluate whether the response by the LMM is informative,
and whether hallucination exists in it, based on the comparison between the LMM's response and the
factual information provided in the image contents, question, and the standard human-generated answer
below.

Please note that the standard human-generated answer may only contain factual information but may not
give a detailed analysis. Also, the standard human-generated answer may not be completely comprehensive
 in describing all the objects and their attributes, so please be a bit more cautious during
evalutation. LMM's detailed analysis or reasoning should be encouraged.

To evaluate the LMM responses, you must rate the response by choosing from the following options:
- Rating: 6, very informative with good analysis or reasoning, no hallucination
- Rating: 5, very informative, no hallucination
- Rating: 4, somewhat informative, no hallucination
- Rating: 3, not informative, no hallucination
- Rating: 2, very informative, with hallucination
- Rating: 1, somewhat informative, with hallucination
- Rating: 0, not informative, with hallucination

Just answer a number in range [0, 6], nothing else.
```

Listing 1. System prompt template for hallucination annotation with DeepSeek-V3

Table 4. Training hyperparameters of different stages.

| Configuration | Classification | Iteration 1 | Iteration 2 |
|---|---|---|---|
| Global batch size | 24 | 24 | 32 |
| Peak learning rate | 1e-4 | 1e-5 | 2e-6 |
| Epochs | 3 | 1 | 5 |
| LoRA rank | | 128 | |
| LoRA $\alpha$ | | 256 | |
| LoRA dropout | | 0.05 | |
| $\beta_1$ | | 0.9 | |
| $\beta_2$ | | 0.999 | |
| $\epsilon$ | | 1e-6 | |
| Optimizer | | AdamW | |
| Learning rate schedule | | cosine decay | |
| Weight decay | | 0.0 | |
| Warmup ratio | | 0.05 | |

To transform image content into text while controlling information loss for DeepSeek-V3 annotation, we extract key objects through COCO labels as image content, since POVID images originate from COCO 2014 [18]. Apart from the system prompt, our input to DeepSeek-V3 includes the following content:

$$Input = \text{### \textbf{Image Contents}} \ \backslash n \ \{image\_content\} \ \backslash n \backslash n$$
$$\text{### \textbf{Question}} \ \backslash n \ \{question\} \ \backslash n \backslash n$$
$$\text{### \textbf{Standard Human-Generated Answer}} \ \backslash n \ \{gt\_answer\} \ \backslash n \backslash n$$
$$\text{### \textbf{LMM Response to Evaluate}} \ \backslash n \ \{model\_answer\}$$

This process generates 8.4K binary classification training samples, with our trained classifier achieving 90% consistency with DeepSeek-V3 judgments on the held-out validation set. **It is worth noting that although we obtained fine-grained score annotations during the labeling stage, we did not directly use these scores.** Instead, we mapped samples with scores from 0 to 2 as hallucinated, and those with scores from 3 to 6 as non-hallucinated. The purpose of fine-grained scoring during annotation was solely to ensure the interpretability and reliability of the labels.

### C.2. Implementation Details

All models are trained using LoRA, with uniform settings of LoRA rank=128, LoRA alpha=256, and LoRA dropout=0.05. For multimodal models, we freeze the vision encoder and fine-tune only the intermediate projection layer and the subsequent language model. The optimizer is consistently set as the Adam optimizer with warmup, using default parameters ($\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 1e - 6, weight\_decay = 0.0, warmup\_ratio = 0.05$), paired with a cosine learning rate schedule. Training for all 7B-sized models utilize DeepSpeed ZeRO-2, while training for the 13B models employed DeepSpeed ZeRO-3.

More specific hyperparameter settings are provided in Table 4. When fine-tuning Qwen2-VL-7B-Instruct as a hallucination classifier, we set the global batch size to 24 and the initial learning rate to 1e-4, training for a total of 3 epochs. For preference optimization, we perform two iterations of training. In the first iteration, we use off-policy data, with a global batch size of 24, an initial learning rate of 1e-5, and train for 1 epoch. We set the DPO coefficient $\beta = 0.5$, and incorporate the NLL loss into the objective as a regularization term, with a weight of 0.2. Adding the NLL loss helps the model better capture the detailed linguistic style of ground truth answers, encouraging the generation of longer responses.
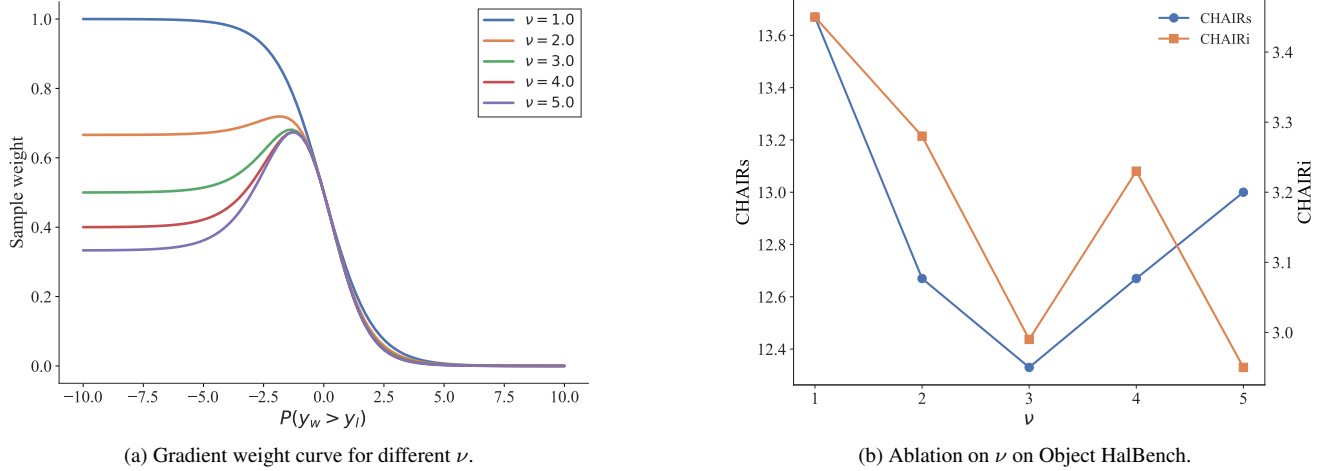
(a) Gradient weight curve for different $\nu$.

(b) Ablation on $\nu$ on Object HalBench.

Figure 5. Ablation study on parameter $\nu$.

The larger $\beta$ further strengthens the KL divergence constraint, contributing to training stability. In the second iteration, we optimize the objective defined in Equation (8), with a global batch size of 32, an initial learning rate of 2e-6, and a total of 5 training epochs. The DPO coefficient is set to $\beta = 0.1$ in this stage.

## D. Additional Results

### D.1. Ablation study on parameter $\nu$

We present additional experiment results examining key parameters in our framework. We mainly present results concerning the dynamic weight model parameters specified in Equation (7). First, we systematically vary $\nu$ to investigate its impact on the gradient difference term $(\nabla_{\theta,y_w} - \nabla_{\theta,y_l})$ in Equation (8). Notably, when $\nu = 1$, our loss function reduces to the standard DPO. Then we conduct ablation on different values of parameter $\nu$ on Object HalBench. The results are shown in Figure 5.

As shown in Figure 5a, the original DPO objective (with $\nu = 1.0$) assigns gradient weights in a sigmoid-like manner across different samples, consistent with the form of the $\sigma(\hat{r}(x, y_l) - \hat{r}(x, y_w))$ term. In contrast, our probability model places greater emphasis on samples near the decision boundary $(P(y_w \succ y_l) \approx 0)$, while assigning lower weights to samples with large negative margins $(P(y_w \succ y_l) \ll 0)$, as such instances are likely to be potentially noisy samples. Figure 5b demonstrates that the dynamic weighting achieves optimal performance at $\nu = 3$, leading us to adopt this value throughout our experiments.

### D.2. General Benchmark Evaluations

To demonstrate that our method effectively reduces model hallucination without compromising general capabilities, we evaluate various hallucination mitigation approaches on both MMBench-EN and MMBench-CN, as shown in the Table 5. The results indicate that our method outperforms baseline models not only in hallucination-related metrics but also in general visual question answering benchmarks. Compared to other algorithms, our method also achieves leading average rankings on both MMBench-EN and MMBench-CN.

## E. Algorithm

We present our complete algorithm in Algorithm 1.

Table 5. Comparison of hallucination mitigation approaches on MMBench-EN and MMBench-CN

| Model Size | Algorithm | Avg. Score ↑ | | Avg. Ranking ↓ |
| --- | --- | --- | --- | --- |
| | | MMBench-EN | MMBench-CN | |
| 7B | LLaVA-Instruct-1.5 [20, 21] | 64.37 | 58.76 | 4.25 |
| | LLaVA-RLHF [35] | 51.40 | 39.52 | 7.0 |
| | HA-DPO [52] | 64.54 | 58.76 | 3.25 |
| | POVID [53] | 64.46 | **60.82** | 2.5 |
| | RLAIF-V [49] | 62.84 | 57.90 | 6.0 |
| | OPA-DPO [47] | **65.73** | 58.42 | 3.0 |
| | **Ours** | <u>65.48</u> | <u>59.36</u> | **2.0** |
| 13B | LLaVA-Instruct-1.5 | <u>67.77</u> | <u>63.75</u> | **1.5** |
| | LLaVA-RLHF | 60.10 | 52.66 | 4 |
| | OPA-DPO | 67.43 | 62.97 | 3 |
| | **Ours** | **69.13** | <u>63.49</u> | **1.5** |

---

**Algorithm 1:** Robust Iterative Alignment

---

**Input:** Classifier $\mathcal{H}$, collected dataset $\mathcal{D} = \{(x, y^*)^i\}_{i=1}^N$, number of iterations $T$, number of generations $K$, batch size $B$, parameter $\nu$ for RK model, learning rate $\eta$.

1 **Initialize**: policy $\pi_{\theta_0}$, preference data set $\mathcal{D}_{\text{pref}} = \emptyset$;
    // Iterative DPO
2 **for** *t=1* **to** *T* **do**
    // Stage 1: Preference Data Construction
3     **for** *i=1* **to** *N* **do**
4         **for** *j=1* **to** *K* **do**
5             generate response $y_j \sim \pi_{\theta_{t-1}}(\cdot \mid x_i)$ for $x_i$ in $\mathcal{D}$ ;   // generate K responses for each prompt
6             Calculate the probability $P(h = 1 \mid x_i, y_j)$ through the hallucination classifier $\mathcal{H}(x_i, y_i^*, y_j)$.
7         **end**
8         Rank hallucination probablities $P(h = 1 \mid \cdot)$ for set $\{x_i, y_j\}_{j=1}^K$;
9         Let the response with highest hallucination probability $P_{\max}$ be $y_l$;
10        Let the response with lowest hallucination probability $P_{\min}$ be $y_w$;
11        **if** $P_{\min} < 0.5$ *and* $P_{\max} \geq 0.5$ **then**
12            $(x_i, y_w, y_l) \to \mathcal{D}_{\text{pref}}$.
13        **end**
14     **end**
    // Stage 2: Robust DPO Training
15     **for** *each epoch* **do**
16         Sample mini-batch $\mathcal{D}_m = \{(x, y_w, y_l)^m\}_{m=1}^B$ from $\mathcal{D}_{\text{pref}}$;
17         Predict the probabilities $\pi_{\theta_t}(y_w \mid x)$ and $\pi_{\theta_t}(y_l \mid x)$ for $(x, y_w, y_l)$ in $\mathcal{D}_m$ using the policy model;
18         Predict the probabilities $\pi_{\theta_{t-1}}(y_w \mid x)$ and $\pi_{\theta_{t-1}}(y_l \mid x)$ for $(x, y_w, y_l)$ in $\mathcal{D}_m$ using the reference model;
19         Calculate the implicit reward $\hat{r}_w = \beta \log \frac{\pi_{\theta_t}(y_w|x)}{\pi_{\theta_{t-1}}(y_w|x)}$, $\hat{r}_l = \beta \log \frac{\pi_{\theta_t}(y_l|x)}{\pi_{\theta_{t-1}}(y_l|x)}$;
20         Calculate pair-wise loss $\ell_{\text{pair}} = \log \sigma(\hat{r}_w - \hat{r}_l)$;
21         Calculate sample weight $\gamma(x, y_w, y_l) = p(y_w \sim y_l \mid x) + \frac{2}{\nu+1}$ ;   // Equation (7)
22         $\theta \leftarrow \theta + \nabla_\theta \mathbb{E}_{(x,y_w,y_l) \sim \mathcal{D}_m} [\text{sg}(\gamma(x, y_w, y_l)) \cdot \ell_{\text{pair}}]$ ;   // Equation (8)
23     **end**
24     $\mathcal{D}_{\text{pref}} = \emptyset$.
25 **end**
**Output:** $\pi_\theta$

Our algorithm implements an iterative model fine-tuning loop that progressively enhances output quality and reduces hallucination through three key phases per iteration. First, the generation phase produces multiple responses per prompt using temperature-controlled sampling to ensure diversity. The subsequent filtering phase applies our adaptive hallucination classifier to exclude hallucinated responses from preference training data. Following each iteration's data collection, we employ a robust reweighting mechanism that dynamically balances reward margin significance to prioritize uncertain boundary samples. This robust reweighting mechanism ensures stable fine-tuning against potential classifier annotation noise.