

RLHF-V: Towards Trustworthy MLLMs via Behavior Alignment from Fine-grained Correctional Human Feedback

Tianyu Yu¹ Yuan Yao^{2*} Haoye Zhang¹ Taiwen He¹ Yifeng Han¹
Ganqu Cui¹ Jinyi Hu¹ Zhiyuan Liu^{1*} Hai-Tao Zheng^{1*} Maosong Sun¹ Tat-Seng Chua²

¹Tsinghua University ²National University of Singapore

yiranytianyu@gmail.com yaoyuanthu@gmail.com

<https://rlhf-v.github.io>

Abstract

Multimodal Large Language Models (MLLMs) have recently demonstrated impressive capabilities in multimodal understanding, reasoning, and interaction. However, existing MLLMs prevalently suffer from serious hallucination problems, generating text that is not factually grounded in associated images. The problem makes existing MLLMs untrustworthy and thus impractical in real-world (especially high-stakes) applications. To address the challenge, we present RLHF-V, which enhances MLLM trustworthiness via behavior alignment from fine-grained correctional human feedback. Specifically, RLHF-V collects human preference in the form of segment-level corrections on hallucinations, and performs dense direct preference optimization over the human feedback. Comprehensive experiments on five benchmarks in both automatic and human evaluation show that, RLHF-V can enable substantially more trustworthy MLLM behaviors with promising data and computation efficiency. Remarkably, using 1.4k annotated data samples, RLHF-V significantly reduces the hallucination rate of the base MLLM by 34.8%, outperforming the concurrent LLaVA-RLHF trained on 10k annotated data. The final model achieves state-of-the-art performance in trustworthiness among open-source MLLMs, and shows better robustness than GPT-4V in preventing hallucinations aroused from over-generalization.

1. Introduction

The recent success of Multimodal Large Language Models (MLLMs) marks a significant milestone in AI research [3, 5, 13, 15, 25, 33, 35, 52, 61]. By connecting visual signals and Large Language Models (LLMs), MLLMs show unprece-

ded capabilities in multimodal understanding, reasoning, and interaction [34, 35, 55]. The models are typically pre-trained on large-scale image-text data to learn the foundational multimodal knowledge and capabilities [3, 5, 15, 25]. To steer the model behavior, most MLLMs are further fine-tuned with instruction tuning (also known as supervised fine-tuning), which supervises models to clone the behavior from demonstration data, enabling MLLMs to engage users in various open-world tasks [5, 13, 31, 33, 58].

However, current MLLM behaviors are not well aligned with human preferences. A glaring issue is their tendency to produce *hallucinations* — responses that are not factually grounded in the associated images [27, 31, 35, 46]. This typically includes descriptions of non-existing visual contents and errors in descriptions. As shown in Figure 1, current MLLMs can hallucinate about objects, attributes, numbers, positions, actions, etc. Quantitatively, our human evaluation shows that the problem is prevalent among state-of-the-art MLLMs, where even the most advanced GPT-4V [35] contains obvious hallucinations in 45.9% responses. The problem makes existing MLLMs untrustworthy and thus impractical in real-world (especially high-stakes) applications, such as guiding visually impaired individuals [35] or autonomous driving systems [53].

We argue that the problem arises from the lack of positive/negative human feedback in instruction-tuned models, making it challenging to learn the precise behavior boundaries to exclude hallucination. To address the problem, we propose RLHF-V, a novel framework that aligns MLLM behavior by learning from human feedback. A straightforward way is to employ the traditional Reinforcement Learning from Human Feedback (RLHF) method in state-of-the-art LLMs [36, 49], which involves human annotators ranking model responses, and utilizing a reward model to guide the policy LLM learning. However, this approach is fraught with two key challenges: (1) *Annotation ambiguity*. Help-

*Corresponding authors

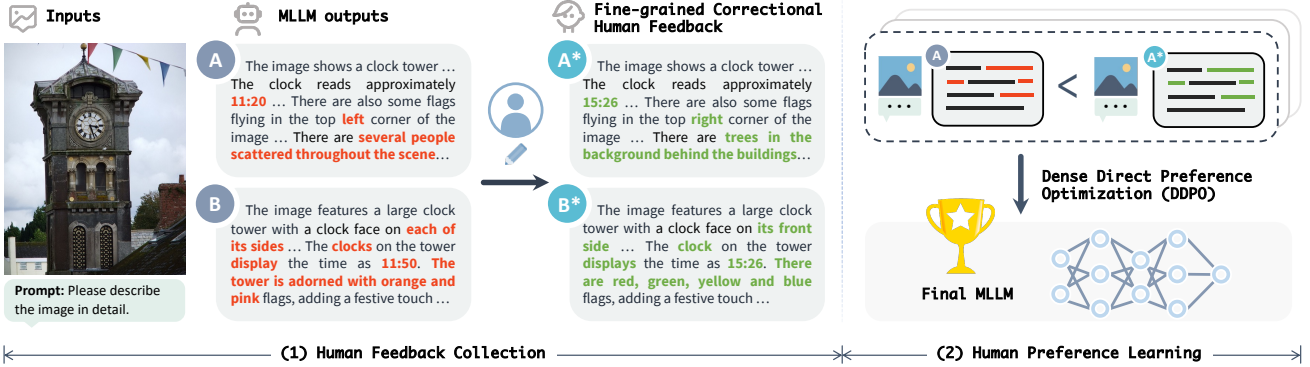


Figure 1. The RLHF-V framework for MLLM behavior alignment from human feedback. (1) Given the input image and prompt, we obtain outputs from MLLMs and collect human feedback in the form of fine-grained segment-level **corrections** on **hallucinations**. (2) During human preference learning, we perform dense direct preference optimization over the fine-grained correctional human feedback.

ful and engaging responses about rich image content are typically long and complex, making it usually non-obvious to decide which response is preferable. As shown in Figure 1 (responses A and B), annotators usually face dilemmas when presenting responses with respective advantages and flaws. Besides, even if labeled with a clear preference, the optimal response remains unknown (e.g., the exact time of the clock). (2) *Learning efficiency*. The coarse-grained ranking feedback makes it difficult to accurately allocate credit to the desirable behaviors. Considering the linguistic complexity and variance of responses, the desirable behavior often requires a large amount of labeled data to learn [12, 37, 46]. Moreover, misallocation of credit to the non-robust bias correlated with the data usually leads to reward hacking and behavior degeneration problems [6, 49].

RLHF-V addresses these challenges by introducing two key innovations: (1) At the data level, we propose to collect human feedback in the form of fine-grained segment-level corrections. As shown in Figure 1, we ask human annotators to directly correct the hallucinated segments from model responses, providing a clear, dense, and fine-grained human preference, as well as optimal responses. This strategy also avoids linguistic variance and non-robust bias, ensuring that the feedback is accurately allocated to the desirable behaviors, thereby enhancing learning efficiency and preventing reward hacking problems. (2) At the method level, we propose dense direct preference optimization (DDPO), a new variant of DPO [40] that addresses the traditional RLHF objective in an equivalent simple and efficient supervised fashion. DDPO directly optimizes the policy model against dense and fine-grained segment-level preference, where the hallucinated segments receive stronger feedback to be factually grounded.

Comprehensive experiments on five benchmarks show that, RLHF-V can substantially enhance the trustworthiness of MLLMs with promising data and computation efficiency.

Using 1.4k preference data, RLHF-V significantly reduces the object hallucination rate of the base MLLM by 34.8%, surpassing the concurrent LLaVA-RLHF [46] trained on 10k preference data. We also show that RLHF-V achieves better robustness than the strong GPT-4V [35] in preventing hallucinations aroused from over-generalization.

The contribution of this work can be summarized as threefold: (1) We present RLHF-V, a novel framework that aligns MLLM behavior through fine-grained correctional human feedback. (2) We collect high-quality human preference data to provide human-aligned learning signals for MLLMs. (3) We conduct comprehensive experiments to demonstrate the effectiveness of the proposed framework, achieving state-of-the-art performance in trustworthiness among open-source MLLMs. All the code, data, and model weights are open-sourced to facilitate future research.

2. Human Preference Collection

The goal of human preference data is to distinguish human-preferred high-quality responses from inferior ones, providing human-aligned learning signals to steer the MLLM behaviors. We first provide an analysis of underlying factors of human preference data, based on which we motivate the human preference collection procedure of RLHF-V.

Human Preference Data: Underlying Factors and Challenges. Given the input x (including the image and the prompt), denote the difference between a preferred output y_w and an inferior output y_l as Y . The difference Y can be essentially decomposed into three factors:

$$Y = Y_p + Y_s + Y_n, \quad (1)$$

where Y_p is the truly preferred behavior such as being trustworthy and helpful, Y_s denotes the shallow non-robust bias correlated with the data but unrelated to human judgment (e.g., y_w contains more usage of specific words), and Y_n

is the random noise factor denoting the linguistic variance of natural language (e.g., different ways of expressing the same meaning). Y_p is the factor we want to learn from the difference Y , while fitting to Y_s can lead to reward hacking problems and thus should be avoided. The linguistic variance Y_n does not bias the preference learning but makes the learning more difficult, demanding more labeled data to learn to the preferred factor Y_p , and thus should also be avoided if possible.

The common RLHF practices in LLMs collect human preference Y in the form of ranking labels, indicating the overall relative quality of responses [36, 37, 49]. According to the above analysis, the practice faces several key challenges: (1) *Annotation ambiguity*. It can be non-obvious to annotate which response is superior using an overall ranking label due to the fine-grained nature of Y_p , especially for complex responses. As shown in Figure 1, annotators usually cannot agree on assigning an overall ranking to different responses with respective advantages and flaws. We observe the issue leads to unsatisfactory annotation quality of existing RLHF data. Moreover, even if labeled with a clear preference, the optimal responses for the questions typically remain unknown. (2) *Learning efficiency*. During reinforcement learning, it can be challenging and data-demanding to precisely allocate the sparse and coarse-grained credit from Y through the linguistic variance Y_n to the preferred behavior Y_p . Misallocation to the non-robust bias factor Y_s will lead models to collapse to exploit trivial rewards [6, 49].

Fine-grained Correctional Human Preference Collection. To address the challenges, we propose to collect fine-grained human preferences in the form of segment-level corrections. As shown in Figure 1, given a flawed output y_l from MLLMs, we ask human annotators to directly correct the hallucinated segments, resulting in a factually optimal output y_w . The annotation simultaneously yields a segment-level incremental preference pair (y_w, y_l) . The simple procedure effectively addresses the challenges: (1) The annotation of incremental correction in segments is clearer and more operable for human labelers. (2) The dense and fine-grained feedback is directly allocated to the preferred behavior Y_p , excluding the linguistic variance Y_n and the non-robust bias Y_s , therefore improving learning efficiency and preventing reward hacking problems. In experiments, we find that the procedure greatly improves the annotation quality and data efficiency, enabling our model to surpass concurrent models trained on an order of magnitude more labeled preference data (see Section 4.3).

In practice, we obtain a total of 1.4k prompts as input from existing instruction tuning dataset [58] and image description prompts generated by GPT-4, and get the responses from Muffin [58] for human annotation. The responses after annotation contain 64.4 words and 2.65 corrected segments on average. We observe that the correc-

tions are diverse in hallucination types, including objects (41.2%), positions (20.3%), numbers (16.5%), attributes (10.0%), actions (5.3%) and miscellaneous types (6.8%).

3. Method

We introduce the RLHF-V approach that learns the fine-grained correctional human feedback by dense direct preference optimization. In addition, we also mitigate existing sources of hallucination in MLLM training by addressing the vision-language mismatch problem.

3.1. Dense Direct Preference Optimization

To leverage the dense and fine-grained human feedback, we present DDPO, a new variant of direct preference optimization [40] for directly optimizing the MLLM policy against dense human preference. The prevalent RLHF approaches involve fitting a reward model on the preference data, and then training the critique, policy and value models to maximize the reward without deviating too far from the reference model [12, 37, 49]. This procedure requires training multiple LLMs with extensive sampling and training, which suffers from complex procedures and high computation cost.

Direct Preference Optimization (DPO) [40] solves this reinforcement learning objective in a simpler equivalent supervised fashion. Here we briefly introduce the DPO method, and refer readers to the original paper for more details. The key observation of DPO is that the reward function $r(x, y)$ can be analytically expressed by its optimal policy model $\pi_*(y|x)$ and reference model $\pi_{\text{ref}}(y|x)$, and therefore we can directly optimize the policy model under proper forms on the preference data. Specifically, the reward model $r(x, y)$ can be represented as:

$$r(x, y) = \beta \log \frac{\pi_*(y|x)}{\pi_{\text{ref}}(y|x)} + \beta \log Z(x), \quad (2)$$

where β is a constant and $Z(x)$ is the partition function. Based on this observation, the policy model can be directly optimized on the human feedback data:

$$\begin{aligned} \mathcal{L} &= -\mathbb{E}_{(x, y_w, y_l)} [\log \sigma(r(x, y_w) - r(x, y_l))] \\ &= -\mathbb{E}_{(x, y_w, y_l)} [\log \sigma(\beta \log \frac{\pi_*(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_*(y_l|x)}{\pi_{\text{ref}}(y_l|x)})], \end{aligned} \quad (3)$$

where the reference model $\pi_{\text{ref}}(y|x)$ is usually implemented by an instruction-tuned base model we want to improve, and is kept fixed during DPO training. Only the policy model $\pi_*(y|x)$ is updated. We note that DPO is more simple, efficient and stable in aligning MLLM behaviors compared with traditional RLHF approaches.

Leveraging dense and fine-grained segment-level feedback essentially requires the model to evaluate the reward of segment-level actions. However, DPO is designed for learning preference in the form of overall response ranking

labels. Specifically, the action score of DPO is given by the likelihood of the holistic response in practice, where different segments are equally treated:

$$\log \pi(y|x) = \sum_{y_i \in y} \log p(y_i|x, y_{<i}), \quad (4)$$

where y_i is the i -th token of the response y . We argue that compared with unchanged segments y_u , corrected segments y_c more directly reveal human judgment in hallucination, and thus should contribute more to the overall action evaluation. Therefore, we propose to score the response as a weighted aggregation of the fine-grained segments:¹

$$\log \pi(y|x) = \frac{1}{N} \left[\sum_{y_i \in y_u} \log p(y_i|x, y_{<i}) + \gamma \sum_{y_i \in y_c} \log p(y_i|x, y_{<i}) \right], \quad (5)$$

where $\gamma > 1$ is a weighting hyperparameter, and larger γ means more contribution from the corrected segments. $N = |y_u| + \gamma|y_c|$ is a normalizing factor, preventing longer responses from getting higher scores. In this way, corrected segments are highlighted to receive stronger human preference feedback to be factually grounded. In experiments, we find that DDPO can better exploit the fine-grained human feedback, leading to more trustworthy responses.

3.2. Mitigating Hallucination from VL Mismatch

DDPO reduces hallucination by learning from human feedback. From another cause-and-effect view, we examine the mainstream MLLM training paradigm, and identify sources of hallucinations in training MLLMs. Based on the observations, we motivate a more trustworthy training recipe.

In general, current MLLMs learn multimodal capabilities in a supervised learning paradigm, where the model outputs are supervised against the ground-truth text associated with the image. In such a paradigm, hallucinations can be introduced by mismatches between images and text data. In practice, the mismatch can come from: (1) low-quality text in pre-training and instruction tuning data, and (2) careless image augmentation during training. We specify the issues and solutions in the following.

Addressing Low-quality Text Influence. Current pre-training data of MLLMs are automatically crawled from the Web [8, 9, 42], which inevitably suffers from severe noise in the text even after extensive post-processing. Supervising MLLMs against such data is essentially teaching them to hallucinate (e.g., describing elements not present in the image, or producing inconsistent descriptions with the image). Similarly, most existing visual instruction tuning datasets are generated by ChatGPT/GPT-4 according to intermediate text annotations [31, 33, 58], which inevitably introduces hallucination into instruction data. While it can be difficult

¹For denotation simplicity, without confusion we also use y_u and y_c to denote the set of tokens in unchanged and corrected segments respectively.

to repair existing pre-training and instruction-tuning data, we find that the influence can be countered by simply post-training MLLMs on high-quality visual question-answering datasets. Intuitively, human-labeled datasets can provide accurate learning signals to calibrate model behaviors from hallucinations, and also enhance instruction-following capabilities. In our experiments, we find that simply fine-tuning the model on VQAv2 [17] can significantly reduce the hallucination rate (see Section 4.3).

Mitigating Untrustworthy Image Augmentation. The vision-language mismatch can also come from the image domain. Data augmentation is widely adopted to improve the data diversity and model robustness in various multimodal models [13, 25, 39, 51, 58]. However, we note that such augmentation must be performed with care in training MLLMs. The key problem is that some image augmentation operations can significantly change the semantics of images, which may make the augmented image inconsistent with the associated text. For example, during augmentation, random cropping can make the objects mentioned in the text absent from the image. This can make the model describe non-existing objects, with wrong numbers, and in wrong positions. In our model training, we exclude image cropping in data augmentation, which improves the trustworthiness of MLLMs (see Section 4.3).

4. Experiments

In this section, we empirically investigate the effectiveness of RLHF-V in aligning MLLM behaviors. In addition to evaluating the trustworthiness and helpfulness of conversation, we also analyze the data efficiency and scalability as well as the robustness. We refer readers to the appendix for more details on benchmarks, baselines and results.

4.1. Experimental Settings

We first introduce the experimental settings, including evaluation, baselines, and implementation details.

Evaluation. We evaluate the models from two perspectives, including trustworthiness reflecting the hallucination degree, and helpfulness reflecting the general interaction quality. Similar to [46], we find binary classification evaluation (i.e., answering yes/no) [16, 27] cannot well reflect the MLLM behaviors in open-ended long-form interactions. We thus adopt benchmarks that directly evaluate the long-form responses, which are more closely related to the practical usage scenarios of MLLMs. For trustworthiness, we perform evaluation on three benchmarks:

(1) **Object HalBench** [41] is a widely adopted benchmark for assessing object hallucination in detailed image descriptions. It compares the objects in the model output with object labels exhaustively annotated for COCO images [29] to detect object hallucination. To improve the

Model	Object HalBench ↓		MHumanEval ↓				MMHal-Bench		LLaVA Bench			VQAv2
	Resp.	Mention	Object	Position	Number	All	Info.	Resp. ↓	Conv.	Detail	Comp.	testdev
LLaVA [33]	63.0	29.5	46.6	21.2	19.9	80.8	31.9	70.8	85.4	74.3	96.3	-
Muffin [58]	50.5	24.5	33.6	16.4	26.0	74.7	33.4	68.8	89.3	79.7	<u>97.7</u>	-
LRV [31]	32.3	22.3	43.2	<u>11.6</u>	19.2	82.9	22.2	78.1	61.7	47.3	55.0	-
LLaVA-RLHF [46]	38.1	18.9	37.7	17.8	18.5	72.6	<u>39.9</u>	65.6	93.8	74.3	111.4	-
InstructBLIP [13]	<u>25.9</u>	<u>14.3</u>	<u>30.8</u>	15.1	17.1	63.7	29.5	<u>64.4</u>	83.2	67.6	90.6	-
Qwen-VL-Chat [5]	43.8	20.0	34.9	16.4	<u>15.8</u>	<u>61.0</u>	38.5	52.1	81.9	<u>77.1</u>	92.3	<u>79.5</u>
LLaVA 1.5 [32]	46.3	22.6	<u>30.8</u>	17.8	17.1	<u>61.0</u>	39.2	52.1	81.6	75.5	95.2	80.0
RLHF-V	12.2	7.5	21.9	7.5	14.4	55.5	40.0	52.1	<u>93.1</u>	75.3	91.6	80.0
GPT-4V [35]	13.6	7.3	22.6	12.3	11.0	45.9	47.6	31.3	96.0	102.5	106.7	77.2*

Table 1. Main experimental results on hallucination. We report hallucination rates in different granularities, including response-level (Resp.) and mention-level (Mention), and response-level hallucination rates in different types. We also show scores on informativeness (Info.), multimodal conversation (Conv.), detailed description (Detail), and complex reasoning (Comp.). * denotes zero-shot results on VQAv2.² The best and second best open-source results are shown in **bold** and underlined respectively.

evaluation stability, we augment the benchmark with 8 diverse prompts for detailed image descriptions. We report the response-level hallucination rate (i.e., the percentage of responses that have hallucinations), as well as the mention-level hallucination rate (i.e., the percentage of hallucinated object mentions among all object mentions).

(2) **MMHal-Bench** [46] evaluates hallucinations and response informativeness. It employs GPT-4 to compare model output with human response and several object labels to decide the scores. In experiments, we find that GPT-4 cannot reliably detect hallucinations due to the incompleteness of MMHal-Bench text annotations. We therefore only report the informativeness score from GPT-4, and assess response-level hallucination rate by human evaluation.

(3) **MHumanEval**. The above evaluations are either limited to common object hallucination or dominated by short-form question answering (i.e., questions that can be sufficiently answered by a few words). To provide a more reliable and comprehensive evaluation over diverse hallucination types, we present MHumanEval benchmark, which covers both long-form image descriptions, and short-form questions. The benchmark contains 146 samples collected from Object HalBench (50) and MMHal-Bench (96). Given model responses, we ask human annotators to label the hallucinated segments and hallucination types of the segments, including objects, positions, numbers and others. We report the response-level hallucination rate on these types.

For helpfulness, we adopt two benchmarks: (1) **LLaVA Bench** [33] is a widely adopted benchmark for assessing multimodal conversation, detailed description and complex reasoning capabilities. It scores model output against reference response via GPT-4. (2) **VQAv2** [17] is a dataset for short-form visual question answering. To facilitate replication and future research, we provide all the model responses

and evaluation results in the supplementary materials.

Baselines. We compare our model with state-of-the-art baselines. (1) **General baselines.** We adopt Qwen-VL-Chat [5], LLaVA [33], LLaVA 1.5 [32], Muffin [58], and InstructBLIP [13] as representative general baselines. These models are mostly pre-trained on large-scale multimodal data, and fine-tuned on high-quality instruction data, achieving strong performance across various multimodal tasks. (2) **Baselines tailored for hallucination problems.** LRV [31] is fine-tuned on 400k instruction data generated by GPT-4, and mitigates hallucination by limiting the response length. The concurrent LLaVA-RLHF [46] employs the strong 13B Vicuna v1.5 [60] (fine-tuned from LLaMA-2 [49]) as LLM backbone. It trains the reward model on 10k human-labeled preference data, and performs proximal policy optimization [43] on 72k factually augmented data. (3) **Commercial Baseline.** We also include GPT-4V [35] as a strong reference, evaluating the gap between the open-source models and state-of-the-art commercial models.

Implementation Details. We implement the RLHF-V framework based on Muffin [58]. The model uses BEiT-3 [51] as the visual module, and 13B Vicuna v1.0 [11] (fine-tuned from LLaMA [48]) as the LLM backbone. The hyperparameter β is 0.5, and the weighting coefficient γ is 5. We train the model with DDPO for 7 epochs, with image resolution 448, learning rate $5e-7$ and batch size 32. The training of RLHF-V is computationally efficient, which takes less than 1 hour on 8 A100 GPUs in total.

4.2. Main Results

The main experimental results are reported in Table 1, from which we observe that: (1) RLHF-V achieves state-of-

²Due to limited instruction-following capability, most MLLMs need to be specifically fine-tuned to produce short-form VQA answers, and therefore cannot achieve reasonable zero-shot performance on VQAv2.

Model	Living Room			Kitchen			Bathroom			Street			$\overline{\Delta}$
	book, person, bed chair, couch, remote			bottle, bowl, cup person, chair, knife			toilet, sink, bottle toothbrush, person, cup			person, car, motorcycle traffic light, handbag, truck			
	H _a	H _s	Δ	H _a	H _s	Δ	H _a	H _s	Δ	H _a	H _s	Δ	
LLaVA-1.5 [32]	25.2	41.8	+16.6	18.9	23.9	+5.0	22.4	30.4	+8.0	20.6	28.0	+7.4	+9.2
LLaVA-RLHF [46]	23.7	34.5	+10.8	13.1	17.4	+4.3	18.2	19.5	+1.4	18.3	22.7	+4.4	+5.2
QWEN-VL [5]	24.5	34.5	+10.0	16.4	20.8	+4.4	21.6	17.5	-4.1	22.5	32.0	+9.5	+5.0
RLHF-V	5.5	8.0	+2.5	3.8	5.9	+2.1	4.1	4.0	-0.1	2.3	4.6	+2.3	+1.7
GPT-4V [35]	8.2	19.4	+11.2	4.6	5.7	+1.1	5.9	13.3	+7.5	4.2	4.6	+0.4	+5.0

Table 2. Experimental results of hallucination from over-generalization on Object HalBench. For each scene, we report the hallucination rate of the top 10 frequent objects on average on the full benchmark (H_a) and under the scene (H_s). Top 6 frequent objects are listed for each scene for brevity. Δ : hallucination rate difference, $\bar{\Delta}$: average difference across the scenes.

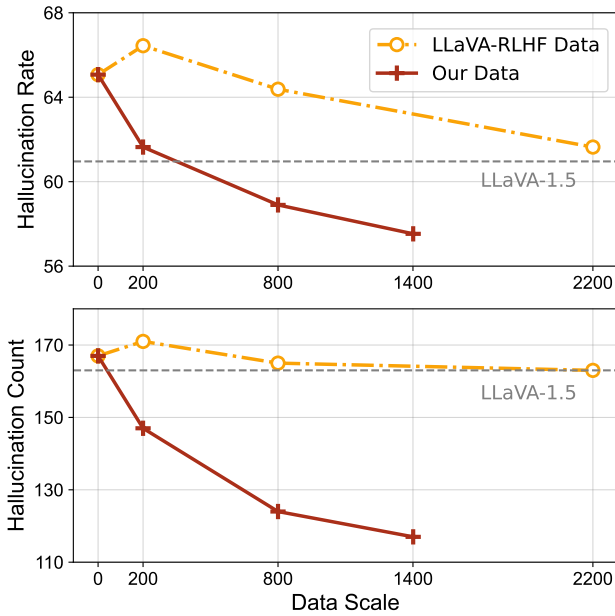


Figure 2. Hallucination rate and number on MHumanEval (all types) with respect to the amount of preference data. We report the results of different models trained on different RLHF data.

the-art performance in trustworthiness among open-source models, outperforming strong general models and models tailored for hallucination. The framework significantly reduces the hallucination rate of the base model Muffin by 75.8% relative points for common objects on Object HalBench, and by 34.8% for overall objects on MHumanEval. The improvement is consistent in different granularities including response-level and mention-level hallucinations, and different hallucination types including objects, positions, and numbers. The reduction is more significant on the more challenging long-form answers on Object HalBench and MHumanEval. The results show that RLHF-V can effectively learn from fine-grained correctional human

feedback to enable more trustworthy MLLM behaviors. (2) RLHF-V achieves promising performance in response helpfulness, where the results on MMHalBench, LLaVA Bench and VQAv2 are strong and comparable to the base model. This shows that RLHF-V can enhance the trustworthiness of MLLMs without sacrificing their helpfulness.

4.3. Analysis

In this section, we conduct analyses on the framework considering the following research questions: (1) How does RLHF-V’s performance scale with feedback data amount? (2) What is the advantage of fine-grained correctional preference data over traditional overall ranking data? (3) Can RLHF-V’s data and method be adopted to enhance the trustworthiness of other MLLMs? (4) How does human feedback alleviate hallucinations intuitively?

Scaling feedback data leads to promising results. We report the hallucination rate and numbers of hallucinated segments on MHumanEval under different amounts of feedback data in Figure 2. We observe that the hallucination rate and number of RLHF-V show a significant and rapid decrease as the data amount grows. This shows that fine-grained correctional human feedback provides effective and efficient learning signals for MLLM behavior alignment. Based on this tendency, we expect better performance can be achieved with an increasing amount of feedback data. We leave this for future work.

Fine-grained correctional human feedback enables better learning efficiency. To quantify the advantage of fine-grained correctional human feedback, we replace our data with the 2.2k human preference data on hallucination from LLaVA-RLHF, which gives overall ranking labels following common RLHF practices. From the experimental results in Figure 2, we observe that model equipped with our data shows a more significant and rapid reduction in hallucination rate and number. Notably, using only 200 preference data, our model achieves comparable hallucination rate to the model that uses an order of magnitude more labeled data

Model	MHumanEval↓				MHB↓	VQAv2
	Obj.	Pos.	Num.	All	Resp.	testdev
Muffin [58]	33.6	16.4	26.0	74.7	68.8	-
RLHF-V	21.9	7.5	14.4	55.5	52.1	80.0
w/ vanilla DPO	21.9	11.6	11.6	57.5	54.2	80.0
w/ IT-VQA only	34.3	17.1	17.1	65.1	58.3	80.0
w/ untrust aug.	18.5	13.7	14.4	59.6	54.2	77.1

Table 3. Ablation results on different components. MHB: MMHal-Bench, IT-VQA: instruction tuning on VQAv2, untrust aug.: untrustworthy data augmentation.

from LLaVA-RLHF. The superior data efficiency is due to (1) better data quality since label ambiguity is minimized, and (2) more direct feedback on hallucinated segments, excluding non-robust bias and linguistic variance.

RLHF-V generalizes to enhance other MLLMs. To investigate the generalization capability of the framework, we adopt RLHF-V’s data and approach to align the behavior of LLaVA [33], a representative and widely used MLLM. Experimental results show that RLHF-V effectively reduces the hallucination count of LLaVA by 13.8 relative points, as well as the hallucination rate by 5.9 relative points. The results demonstrate that RLHF-V is applicable across different MLLMs to improve trustworthiness.

RLHF-V reduces hallucination from correlation and over-generalization. LLMs possess rich world knowledge and strong generalization capabilities. Without proper positive/negative human feedback, MLLMs can over-generalize to produce highly correlated and plausible concepts, which leads to hallucinations. For example, a prevalent hallucination case observed across different MLLMs is claiming the presence of *person* as long as they see an image of *street*. To quantify the problem, we select a set of representative scenes $\{\textit{living room}, \textit{kitchen}, \textit{bathroom}, \textit{street}\}$. For each scene, we identify the corresponding images in COCO by lexically matching the captions with the scene name. Then we obtain the top 10 frequent objects in the scene from the COCO object annotations. We compare the response-level hallucination rate for these objects (1) on average across all test samples, and (2) on samples under the target scene. Models prone to over-generalization will expect a significant increase in the hallucination rate (Δ).

From the experimental results in Table 2, we observe that: (1) All models including GPT-4V show a substantial increase in the hallucination rate, which demonstrates the over-generalization hypothesis. (2) RLHF-V exhibits the smallest change in the hallucination rate, which is even more robust than GPT-4V. The reason for the robustness is that RLHF-V provides crucial positive/negative fine-grained correctional human feedback for MLLMs, which helps to learn clear behavior boundaries between reasonable generalizations and over-generalizations. (3) RLHF-

V achieves the lowest hallucination rates for these common objects both on average and especially under common scenes. This makes RLHF-V preferable in practical real-world applications.

Ablation Study. To investigate the contribution of each component, we perform an ablation study. From the experimental results in Table 3, we can observe that: (1) Learning human feedback with vanilla DPO leads to performance degrades, showing the advantage of DDPO in exploiting the fine-grained human preference. (2) Fine-tuning on VQAv2 leads to a significant reduction in hallucination rates compared with the base model. This reveals the value of traditional human-annotated datasets from a new perspective of hallucination mitigation. (3) Including untrustworthy data augmentation (i.e., image cropping) in training hurts the performance on both hallucination and VQAv2. This shows that careless data augmentation can be a double-edged sword in training MLLMs.

Case Study. To provide an intuitive understanding and comparison of different models, we provide qualitative results in Figure 3. We show cases in two representative scenarios: (1) Short-form QA (i.e., questions that can be sufficiently answered in a few words). Our model typically maintains a good balance between helpfulness, engagement and clarity. In comparison, LLaVA-RLHF is usually far more engaging, introducing lengthy extensions however that can be less reasonable or relevant. (2) Long-form QA (i.e., questions that require long text to answer). We observe that MLLMs are significantly more prone to hallucinations in long-form QA, since it typically requires more comprehensive capabilities from multiple perspectives. For example, InstructBLIP and LLaVA-RLHF can confidently describe non-existing objects in a large proportion of their responses, whereas RLHF-V introduces significantly fewer hallucinations while delivering a comparable amount of effective information. We refer readers to the appendix for more qualitative results.

5. Related Work

Multimodal Large Language Models. Recent trends in multimodal learning have witnessed the success of building MLLMs by connecting visual encoders with powerful LLMs [10, 19, 24, 56, 59]. The current MLLM training paradigm typically involves two stages: (1) Pre-training. Models are pretrained on large-scale image-text pairs [5, 13, 25, 52, 58] or interleaved data [3, 4, 19] to learn the semantic mapping between visual and text signals. (2) Instruction Tuning. To enable the model with instruction-following capability, MLLMs are further fine-tuned on visual instruction data, including collections of existing human-annotated datasets [13, 26, 32], and generated data from ChatGPT/GPT-4 [26, 31, 33, 58]. Despite


<p>Short-form QA</p>  <p>At what time is this photo taken?</p>	<p>RLHF-V (This work) The scene appears to be taken during the daytime, as the sunlight can be seen streaming in through the large window in the kitchen. The bright lighting creates a welcoming and inviting atmosphere, making the kitchen feel more spacious and vibrant.</p> <p>InstructBLIP The photo is taken during the daytime.</p> <p>LLaVA-RLHF The photo is taken during the daytime, as indicated by the presence of natural light and the absence of any artificial lighting. The sunlight streaming in through the windows and the overall brightness of the room suggest that the photo is taken during daylight hours. This could be either in the morning or afternoon, depending on the specific time of day and the position of the sun relative to the windows. The presence of various objects in the room, such as bottles, cups, bowls, and a potted plant, further supports the idea that the photo is taken during the daytime. These objects are likely to be in use or placed in their respective locations during the daytime, as people tend to engage in more activities and prepare meals during the day. In conclusion, the photo is taken during the daytime, as evidenced by the natural light, the absence of artificial lighting, and the presence of various objects in the room that are typically used during daylight hours.</p>
<p>Long-form QA</p>  <p>What is this photo about? Please answer in detail.</p>	<p>RLHF-V (This work) The image depicts a small black dog sitting on a wooden floor in a living room, watching a television screen. The television appears to be mounted on the wall, and it's showing an image of another dog.</p> <p>InstructBLIP The image depicts a black dog sitting on the floor in front of a large TV screen. The TV screen is displaying an image of a dog, which appears to be the same breed as the dog in the room. The dog seems to be engrossed in the TV program and is focused on the screen. There are several other dogs visible in the room, but they are not as close to the TV as the main subject. One dog is located near the left side of the TV, while another is positioned closer to the center of the room. A third dog can be seen towards the right side of the image. In addition to the dogs, there are two chairs in the room, one located near the left side of the TV and the other near the center of the room.</p> <p>LLaVA-RLHF The image features a black dog sitting on a hardwood floor in a room, watching a television. The TV is located on the left side of the room, and the dog is positioned near the center of the scene. There are also two people in the room, one on the left side and the other on the right side, possibly watching the television along with the dog. Additionally, there are two remotes visible in the room, one near the center and the other on the right side. A chair can be seen in the room, placed on the right side, and a bowl is located on the left side of the room, possibly for the dog's food or water.</p>

Figure 3. Qualitative results of different models on short-form QA and long-form QA. **Correct answers**, **unreasonable extensions** and **hallucinations** are highlighted in color respectively.

the success, current MLLMs suffer from serious hallucination problems [27, 30, 31, 46]. Notably, even after extensive efforts, GPT-4V has still been found to be prone to hallucinations, making basic factual errors confidently [35]. The problem undermines practical applications of MLLMs especially in high-stakes scenarios, which has recently drawn increasing attention from the community.

Behavior Alignment for LLMs. Aligning language agent behaviors with human preference has emerged as a promising research direction [21, 23]. Pivotal approaches in LLMs include instruction tuning (or supervised fine-tuning) and RLHF [37, 45]. While supervised fine-tuning is suitable for basic behavior alignment [14, 47], due to the mismatch between likelihood maximization objective and human preference, it may introduce or amplify hallucination [36, 37]. Therefore, RLHF is widely accepted for further behavior and preference alignment [7, 12, 36], where proximal policy optimization (PPO) [43] is recognized as the major technique. Later adaptations attempt to stabilize the optimization process [40] and enclose more fine-grained signals [28, 54]. However, RLHF has rarely been explored in MLLMs to align model behaviors with humans.

Reducing Hallucination for MLLMs. Some prelimi-

nary efforts have been made to alleviate hallucination problems in MLLMs. LRV [31] generates instruction data with negative responses, and mitigates hallucination by limiting the response length. However, limiting the response length does not essentially address the problem, and also undermines the response helpfulness. VIGC [50] iteratively refines the instruction data for better instruction tuning. Woodpecker [57] proposes to post-edit hallucinations by merging the output of MLLMs and a more accurate expert VQA model using GPT-3.5. The post-editing procedure involves external tools and LLMs much larger than the target MLLM online in multiple stages, which leads to high inference costs and delays. Gunjal *et al.* [18] distinguishes the inaccurate parts in responses via human annotation, and internally discourages the hallucinated parts by direct preference optimization. However, the positive behaviors for hallucinated parts are unknown, making the human feedback not complete enough to learn the behavior boundary. The concurrent LLaVA-RLHF [46] employs the traditional RLHF approach [37] on MLLMs, and augments the reward model with rich additional text descriptions. It is therefore similarly challenged with label ambiguity, learning efficiency, and complex training. In comparison, RLHF-V

presents the first fine-grained correctional human feedback learning framework for behavior alignment, and systematically addresses different hallucination sources in training MLLMs, achieving strong performance in trustworthiness.

6. Conclusion

Hallucination is a critical problem preventing practical applications of MLLMs in real-world scenarios. In this work, we present RLHF-V, a novel framework that enhances the trustworthiness of MLLMs by behavior alignment from fine-grained correctional human feedback. Comprehensive experimental results show that our model achieves state-of-the-art performance in trustworthiness especially in challenging long-form responses while maintaining strong helpfulness. In this work, we collect correctional feedback from human annotators. In future, with the progress of more trustworthy and capable MLLMs, we will explore collecting accurate preferences from MLLMs, which can facilitate large-scale preference learning for stronger behavior alignment. Besides, we note that the framework of RLHF-V can potentially also help reduce the hallucinations in LLMs, which we will explore in future.

Contributions

The authors' contributions can be outlined as follows:

- In initializing the project, Yuan Yao and Tianyu Yu design the framework to collect correctional human feedback. Tianyu Yu devise the DDPO algorithm. Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun and Tat-Seng Chua offer invaluable guidance in project design.
- In data collection, Taiwen He, Haoye Zhang, Tianyu Yu and Yuan Yao take charge of the annotation process to ensure the data quality.
- In model training and evaluation, Tianyu Yu implements the training framework. Tianyu Yu, Haoye Zhang and Yuan Yao design the evaluation framework. Tianyu Yu and Haoye Zhang implement the evaluation codebase.
- In paper writing, Yuan Yao and Tianyu Yu write the paper. Haoye Zhang, Taiwen He, Yifeng Han, Ganqu Cui, Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun and Tat-Seng Chua offer suggestions to polish the writing.
- For public usability, Tianyu Yu, Yifeng Han, Jinyi Hu and Yuan Yao promote the open-source project.
- Throughout the project, Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun and Tat-Seng Chua provide invaluable guidance and advice.

References

- [1] GPT-4 vision document. <https://platform.openai.com/docs/guides/vision/calculating-costs>. Accessed: 2023-11-20. 12
- [2] John schulman - reinforcement learning from human feedback: Progress and challenges. https://www.youtube.com/watch?v=hhiLw5Q_UFg. Accessed: 2023-11-20. 12
- [3] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *NeurIPS*, 35: 23716–23736, 2022. 1, 7
- [4] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. OpenFlamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023. 7
- [5] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023. 1, 5, 6, 7, 13
- [6] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislaw 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, 3
- [7] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional AI: Harmlessness from AI feedback. *arXiv preprint arXiv:2212.08073*, 2022. 8
- [8] Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. COYO-700M: Image-text pair dataset, 2022. 4
- [9] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, pages 3558–3568, 2021. 4
- [10] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. PaLI: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022. 7
- [11] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing GPT-4 with 90%* ChatGPT quality, 2023. 5
- [12] Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*, 2023. 2, 3, 8
- [13] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. InstructBLIP: Towards general-purpose vision-language models with instruction tuning, 2023. 1, 4, 5, 7, 13
- [14] Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and

- Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*, 2023. 8
- [15] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. PaLM-E: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023. 1
- [16] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, et al. MME: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023. 4
- [17] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*, pages 6904–6913, 2017. 4, 5, 13, 15
- [18] Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and preventing hallucinations in large vision language models. *arXiv preprint arXiv:2308.06394*, 2023. 8, 13
- [19] Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, et al. Language is not all you need: Aligning perception with language models. *arXiv preprint arXiv:2302.14045*, 2023. 7
- [20] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, 2021. If you use this software, please cite it as below. 13
- [21] Zachary Kenton, Tom Everitt, Laura Weidinger, Iason Gabriel, Vladimir Mikulik, and Geoffrey Irving. Alignment of language agents. *arXiv preprint arXiv:2103.14659*, 2021. 8
- [22] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International Journal of Computer Vision*, 128(7):1956–1981, 2020. 13
- [23] Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. Scalable agent alignment via reward modeling: a research direction. *arXiv preprint arXiv:1811.07871*, 2018. 8
- [24] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023. 7
- [25] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. 1, 4, 7
- [26] Lei Li, Yuwei Yin, Shicheng Li, Liang Chen, Peiyi Wang, Shuhuai Ren, Mukai Li, Yazheng Yang, Jingjing Xu, Xu Sun, et al. M3IT: A large-scale dataset towards multi-modal multilingual instruction tuning. *arXiv preprint arXiv:2306.04387*, 2023. 7
- [27] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023. 1, 4, 8, 13
- [28] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*, 2023. 8
- [29] 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 *ECCV*, pages 740–755. Springer, 2014. 4, 13
- [30] Fuxiao Liu, Tianrui Guan, Zongxia Li, Lichang Chen, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. Hallusion-Bench: You see what you think? Or you think what you see? An image-context reasoning benchmark challenging for GPT-4V(ision), LLaVA-1.5, and other multi-modality models. *arXiv preprint arXiv:2310.14566*, 2023. 8
- [31] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023. 1, 4, 5, 7, 8, 13
- [32] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023. 5, 6, 7, 13
- [33] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 1, 4, 5, 7, 13, 15
- [34] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. MathVista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023. 1
- [35] OpenAI. GPT-4V(ision) system card. 2023. 1, 2, 5, 6, 8, 12
- [36] OpenAI. GPT-4 technical report, 2023. 1, 3, 8
- [37] 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. *NeurIPS*, 35:27730–27744, 2022. 2, 3, 8
- [38] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of ICCV*, 2015. 13
- [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763. PMLR, 2021. 4
- [40] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*, 2023. 2, 3, 8

- [41] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. In *EMNLP*, pages 4035–4045, 2018. 4, 13
- [42] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. LAION-5B: An open large-scale dataset for training next generation image-text models. *NeurIPS*, 35:25278–25294, 2022. 4
- [43] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 5, 8, 13
- [44] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A benchmark for visual question answering using world knowledge. In *Proceedings of ECCV*, 2022. 13
- [45] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *NeurIPS*, 33:3008–3021, 2020. 8
- [46] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented RLHF. *arXiv preprint arXiv:2309.14525*, 2023. 1, 2, 4, 5, 6, 8, 13
- [47] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. Stanford Alpaca: An instruction-following LLaMA model, 2023. 8
- [48] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 5
- [49] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. LLaMA 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 1, 2, 3, 5
- [50] Bin Wang, Fan Wu, Xiao Han, Jiahui Peng, Huaping Zhong, Pan Zhang, Xiaoyi Dong, Weijia Li, Wei Li, Jiaqi Wang, et al. VIGC: Visual instruction generation and correction. *arXiv preprint arXiv:2308.12714*, 2023. 8
- [51] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. Image as a foreign language: BEiT pretraining for vision and vision-language tasks. In *CVPR*, pages 19175–19186, 2023. 4, 5
- [52] Wei Han Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. CogVLM: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023. 1, 7
- [53] Licheng Wen, Xueming Yang, Daocheng Fu, Xiaofeng Wang, Pinlong Cai, Xin Li, Tao Ma, Yingxuan Li, Linran Xu, Dengke Shang, et al. On the road with GPT-4V (ision): Early explorations of visual-language model on autonomous driving. *arXiv preprint arXiv:2311.05332*, 2023. 1
- [54] Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A Smith, Mari Ostendorf, and Hannaneh Hajishirzi. Fine-grained human feedback gives better rewards for language model training. *arXiv preprint arXiv:2306.01693*, 2023. 8
- [55] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. The dawn of LMMs: Preliminary explorations with GPT-4V(ision). *arXiv preprint arXiv:2309.17421*, 9, 2023. 1
- [56] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mPLUG-Owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023. 7
- [57] Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. Woodpecker: Hallucination correction for multimodal large language models. *arXiv preprint arXiv:2310.16045*, 2023. 8
- [58] Tianyu Yu, Jinyi Hu, Yuan Yao, Haoye Zhang, Yue Zhao, Chongyi Wang, Shan Wang, Yinxv Pan, Jiao Xue, Dahai Li, et al. Reformulating vision-language foundation models and datasets towards universal multimodal assistants. *arXiv preprint arXiv:2310.00653*, 2023. 1, 3, 4, 5, 7, 13
- [59] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. LLaMA-Adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023. 7
- [60] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhonghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging LLM-as-a-judge with MT-Bench and Chatbot Arena, 2023. 5, 13
- [61] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. MiniGPT-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 1

A. Zoom-in Study regarding GPT-4V

We perform a zoom-in study of RLHF-V concerning GPT-4V to provide a better understanding of their behaviors.

A.1. Hallucination Patterns

We conduct a comparative analysis of the responses generated by RLHF-V and GPT-4V, and have the following key observations:

(1) Compared with RLHF-V, GPT-4V tends to describe more details in the images and elaborate more on the interrelations among them. Quantitatively, we utilize ChatGPT to extract all the object mentions in the responses of GPT-4V, and find that the average number per response is 2.1 times larger than RLHF-V. We mainly attribute this to the higher resolution (7.6 times larger than RLHF-V) [1] and the more powerful LLM backbone [35].

(2) GPT-4V’s hallucinations are more concentrated in some responses. In HumanEval, the hallucination rates of GPT-4V on *Object* and *Position* are comparable to RLHF-V. However, in the comprehensive *ALL* metric, the hallucination rate is 17.3% lower than RLHF-V. To better understand the reasons behind this phenomenon, we conduct a thorough analysis of the evaluation results. We observe that different types of hallucinations in GPT-4V are often concentrated in a small subset of responses, while contributing to hallucination rates across multiple subcategories. Quantitatively, we sort the responses of each model by the hallucination count in descending order, and plot the curve of hallucination count ratio vs hallucination response ratio. From the results in Figure 4, we can see that the top 45.6% hallucinated responses of GPT-4V contribute to 75% hallucination counts. In comparison, the top 64.6% hallucinated responses of RLHF-V contribute to 75% hallucinations. We refer readers to Section B for more qualitative results.

A.2. Distillation against GPT-4V

Upon observing GPT-4V’s superior fine-grained image perception and text generation capabilities, an intuitive question is, will it be beneficial to distill GPT-4V capabilities through visual instruction tuning? To this end, we collect 1.2k visual instruction data about long-form image descriptions from GPT-4V. We then use the response generated by GPT-4V to fine-tune our model. We observe that the average number of object mentions in the model response significantly increases by 1.8 times compared with the origin model. However, this can be a double-edged sword: as shown in Table 4, the hallucination rate significantly increases as well.

The results are consistent with the hypothesis of [2]: “If we supervise the model against instruction data that far exceeds its own foundational capabilities, we are essentially teaching the model to hallucinate.” Specifically, our

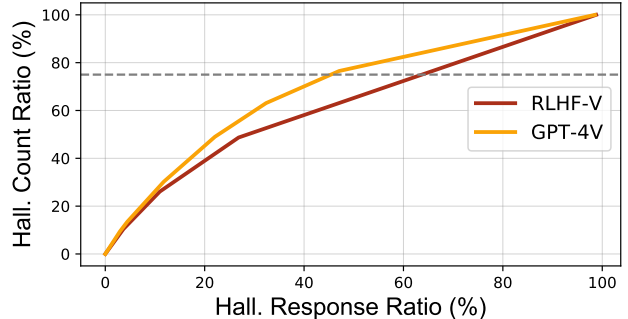


Figure 4. Distribution of hallucination segments over different responses. GPT-4V hallucinations are more concentrated on a smaller subset of the responses. Hall.: Hallucination.

model learns to produce more details and the interrelationship among them through distillation against GPT-4V, while the fundamental capabilities of the model are not enough for this demand. As a result, the hallucination problem is remarkably exacerbated. The results show that visual instruction data (or distillation target) is not the stronger the better, but rather should match the foundational capability of the model.

B. Qualitative Results

We provide more qualitative results in this section to facilitate a more intuitive understanding and comparison of different models. Based on the qualitative results, we have the following observations:

(1) RLHF-V typically exhibits less hallucination in both short-form QA and long-form QA scenarios, compared with open-source models such as LLaVA-RLHF and InstructBLIP, as shown in Figure 5, 6, 7, and 8.

(2) GPT-4V is more descriptive regarding details in images as shown in Figure 6, 8, 9 and 10. For example, in Figure 9, GPT-4V mentions *black dots* across each *tile* while RLHF-V does not describe these details.

(3) RLHF-V is more resistant to the over-generalization problem as shown in Figure 9 and Figure 10. In Figure 9, GPT-4V falsely mentions objects which are highly related to the scene while not shown in the image such as *exhaust*, *hood*, and *bottle*.

C. Implementation Details

We provide more implementation details in this section for better reproducibility. Benefiting from the high efficiency of training, we make all parameters trainable during the training process, which costs merely less than 1 hour on 8 A100 GPUs in total. We empirically find that adopting a longer warm-up (10% training steps) can make the training more stable and consequently apply this setting for all ex-

Model	HumanEval↓				MHB↓
	Obj.	Pos.	Num.	All	Resp.
Muffin [58]	33.6	16.4	26.0	74.7	68.8
RLHF-V	21.9	7.5	14.4	55.5	52.1
w/ GPT-4V distil.	45.2	10.3	20.6	75.3	62.5

Table 4. Experimental results of distillation against GPT-4V. MHB: MMHal-Bench, GPT-4V distil.: instruction-tune the model using responses generated by GPT-4V.

periments in this paper. As for data collection, besides the prompts obtained from [58], we also use image description prompts generated by GPT-4 during the annotation process which are listed in Table 5.

D. Evaluation Details

We introduce more evaluation details, including baseline models and evaluation benchmarks.

D.1. Baselines

We compare with a series of state-of-the-art baselines:

- **LLaVA**: LLaVA [33] constructs 150K multimodal instructions based on the COCO dataset by asking GPT-4 to generate multi-turn dialogues for each image.
- **Muffin**: Muffin [58] propose to reformulate pre-trained vision-language models as the bridge toward large language models. The model is firstly pre-trained on 180M image-text pairs and then fine-tuned on their proposed UniMM-Chat instruction dataset consisting of 1.1M multimodal instructions.
- **LRV**: LRV [31] is fine-tuned on 400K instruction data generated by GPT-4, and mitigates hallucination by limiting the response length.
- **LLaVA-RLHF**: The concurrent LLaVA-RLHF employs the strong 13B Vicuna 1.5 [60] (fine-tuned from LLaMA-2) as LLM backbone. It first trains the model with 122K instructions from VQAv2 [17], A-OKVQA [44] and Flickr30k [38] to improve the foundational capabilities of the model. It then trains the reward model on 10K human-labeled preference data, and performs proximal policy optimization [43] on 72K factually augmented data.
- **InstructBLIP**: InstructBLIP [13] constructs a multimodal instruction tuning dataset based on 26 public datasets by apply pre-defined templates to directly formulate these datasets into a unified format. They also devise a novel instruction-aware Q-Former and train the model on the proposed dataset.
- **Qwen-VL-Chat**: Qwen-VL-Chat [5] utilizes a large ViT with 1.9B parameters initialized from OpenCLIP’s bigG [20] as image encoder. It is pre-trained on 1.4B

image-text pairs and fine-tuned on more than 50M high-quality multimodal instructions.

- **LLaVA 1.5**: LLaVA 1.5 [32] also employs the strong 13B Vicuna 1.5 [60] (fine-tuned from LLaMA-2) as LLM backbone. It is pre-trained on 558K selected image-text pairs and fine-tuned on 665K multimodal instructions with elaborately designed training strategies.

D.2. Benchmarks

We introduce additional details about the benchmarks we used for evaluation.

- **Object HalBench**: Object HalBench [41] is a widely adopted benchmark for assessing object hallucination in detailed image descriptions. To improve the evaluation stability, we augment the benchmark with 8 diverse prompts for detailed image descriptions during evaluation, where 4 instructions are adopted from [18] and the other 4 instructions are generated by GPT-4. We confirm that there is no overlap between the evaluation instructions and the training instructions. Detailed instructions are listed in Table 6. Following [27], we randomly sample 300 images from the validation set of COCO [29] to form the evaluation image set. Regarding metrics, the response-level hallucination rate is the number of responses with object hallucinations divided by the number of responses that introduce COCO objects, while the mention-level hallucination rate is the number of falsely mentioned COCO objects in the generated responses divided by the total number of mentioned COCO objects. During evaluation, we first generate descriptions on images from the benchmark and then leverage ChatGPT to extract the mentioned objects in these responses which are further used to calculate the final scores following. Unlike [41] which detects object mentions by exact-match, we find ChatGPT can perform the extraction with both better precision and recall and consequently apply this setting during evaluation. The full prompt we used to conduct such extraction is shown in Table 7.
- **MMHal-Bench**: MMHal-Bench [46] evaluates hallucinations and response informativeness. It consists of 96 images from the validation and test sets of OpenImages [22]. Each image in this benchmark is annotated with a brand new question and the image-question pairs cover 12 common object meta-categories from COCO.
- **HumanEval**: The above evaluations are either limited to common object hallucination or dominated by short-form question answering (i.e., questions that can be sufficiently answered by a few words). To provide a more reliable and comprehensive evaluation over diverse hallucination types, we present HumanEval benchmark, which covers both long-form image descriptions, and short-form questions. The benchmark contains 146 samples collected from Object HalBench (50) and MMHal-Bench

- Identify and describe each object in the image in detail.
- Describe the key features of the image in great detail.
- What are the main elements in this image? Describe them thoroughly.
- Explain what’s happening in the image with as much detail as possible.
- Detail the image’s components with particular focus on each entity.
- Provide an intricate description of every entity in the image.
- What are the main objects or subjects in the image? Please describe them in detail.
- What is the setting or environment in which the image takes place?
- How do the elements in the image relate to each other in terms of positioning or composition?
- Explain the elements of the image with thorough attention to detail.
- Explain the image’s various components in depth.
- What are the key features you observe in the image?
- Can you point out the details that make this image unique?
- Itemize the elements you identify in the image and describe them thoroughly.
- Convey the specifics of the image with meticulous attention to detail.
- Tell me what catches your eye in the image, and describe those elements in depth.

Table 5. The list of instructions for detailed image description used in training.

- Provide a thorough description of the given image.
- What is this photo about? Please answer in great detail.
- Provide a thorough description of the given picture.
- Explain the narrative or story that the image seems to convey, detailing each part that contributes to it.
- Please provide a detailed description of the image. Describe the visual elements, colors, shapes, textures, and any objects or people present along with the overall mood or atmosphere portrayed in the image.
- Please provide a detailed description of the image, including its visual elements, such as colors, shapes, textures, objects, and people.
- Provide an intricate description of the image, capturing its visual elements, including colors, shapes, textures, objects, and any people present.
- Compose a detailed account of the image, encompassing its visual characteristics, like colors, shapes, textures, objects, and any human subjects, by paying careful attention to the specifics.

Table 6. The list of instructions for Object HalBench.

You are an expert in image objects extraction according to a question answer pair. We asked an examiner to answer a question about a picture.

[Start of Question]

<image> {question}

[End of Question]

[Start of Examiner’s Answer]

{answer}

[End of Examiner’s Answer]

Assume that the answer is correct, please identify all visible objects that are directly shown in the image. Please following the instructions in below:

1. You should only mention objects that are explicitly mentioned in the examiner’s answer.
2. You should only extract the object names without the attributes of the objects.
3. You should not include the properties of the object, like the color, material, etc. as part of the object name in your result.
4. Make your answer precise. Present the results in a JSON list format: [“object_1”, ..., “object_n”].
5. You should return an empty JSON list () if no visible objects can be found.

Table 7. The prompt we used to extract object mentions from image captions with ChatGPT.

(96). Given model responses, we ask human annotators to label the hallucinated segments and hallucination types of the segments, including objects, positions, numbers and others. We report the response-level hallucination rate on these types.

- **LLaVA Bench:** LLaVA Bench [33] is a widely adopted benchmark for assessing multimodal conversation, detailed description and complex reasoning capabilities. It consists of 30 image-question pairs for the aforementioned three capabilities respectively and scores model output against reference response via GPT-4.
- **VQAv2:** VQAv2 [17] is a dataset for short-form visual question answering. The test-dev set of VQAv2 consists of 107K image-question pairs which covers a diverse range of capabilities.



Short-form QA



What is the man doing in the image?

⌘ RLHF-V (This work)

In this picture, the man is **giving a thumbs up gesture**. The thumbs up sign is a universally recognized symbol of approval or agreement, often used to indicate satisfaction with something. It seems like he may be reacting positively to something that is happening around him.

⌘ InstructBLIP

The man in the image is sitting at a table, **talking on his cell phone**.

⌘ LLaVA-RLHF

The man in the image is **talking on his cell phone** while sitting in a crowded room. He is surrounded by other people, **some of whom are also sitting and talking on their phones**. The scene suggests that the man is engaged in a conversation, possibly discussing something important or casually chatting with someone on the other end of the line.

The presence of **multiple people using their cell phones** in the same room indicates that the image may have been taken in a public setting, such as a restaurant, a waiting area, or a social event. It is common for people to use their phones in such settings, as it allows them to stay connected with friends, family, or colleagues, or to attend to work-related matters.

The man's conversation could be a personal one, or it could be related to a professional matter, such as discussing a project or coordinating plans with a colleague. It is also possible that **he is using his phone to access information, such as checking emails, browsing the internet, or using social media**.

In any case, the image captures a moment in time where the man is actively engaged in a conversation on his cell phone, surrounded by **other people who are also using their devices to communicate or attend to their needs**. This scene reflects the ubiquity of cell phones in modern life, and the way they have become an integral part of our daily routines and social interactions.

⌘ GPT4-V

The man in the image appears to be **scratching or touching his head** with one hand while **giving a thumbs-up gesture** with the other hand.

Figure 5. Qualitative results of different models. **Correct answers** and **hallucinations** are highlighted in color respectively.



Long-form QA



What do you think is going on in this snapshot?

⚙️ RLHF-V (This work)

In this image, we see a small gray elephant standing alone in a dirt field. The elephant appears to be the focal point of the photo, with its impressive size and unique color. In the background, there is a building that complements the natural surroundings very well.

⚙️ InstructBLIP

The image depicts an elephant standing in a sandy area near a wall. The elephant is the main focus of the scene, **occupying most of the space in the frame**. There are several other objects scattered around the area, including **a few chairs placed near the elephant**. In addition to the elephant and **chairs**, there are **two tires** visible in the scene. **One tire is positioned closer to the left side of the image**, while the other is located on the right side. **Both tires appear to be partially buried in the sand**, possibly indicating that they have been placed there for the elephant to interact with. Overall, the scene provides a glimpse into the life of an elephant in a controlled environment.

⚙️ LLaVA-RLHF

The image features a large elephant standing in a sandy area, surrounded by rocks and a building. The elephant is positioned near the center of the scene, with its trunk hanging down.

In the background, there are two tires lying on the ground, possibly indicating that the area is used for various activities or training exercises. Additionally, there is **a bird perched on a rock** in the scene, adding a touch of nature to the environment.

⚙️ GPT4-V

The image depicts an elephant in an enclosure, likely at a zoo or a sanctuary. The environment has a concrete wall, a door, and some elements like rocks, a blue tarp, and a tire, possibly used as play or enrichment items for the elephant. The scene suggests that this is a space where the elephant is housed, and the items around might be used to keep it engaged or stimulated.

Figure 6. Qualitative results of different models. **Hallucinations** are highlighted in color.



Long-form QA



What might be the significance or role of the objects you described in ensuring city traffic flows smoothly?

⊗ RLHF-V (This work)

The objects in the image hold significance in ensuring city traffic flows smoothly as they help regulate and control the movement of vehicles and pedestrians. The traffic lights, for instance, provide visual cues to drivers and pedestrians, indicating when it is safe to proceed, stop, or yield. By doing so, they help reduce traffic congestion, prevent accidents, and improve overall safety.

Additionally, the street markings, such as the crosswalk and bike lane, provide clear guidance to drivers, cyclists, and pedestrians by outlining designated areas for each group. These markings contribute to a more organized and safer traffic flow, as they indicate the expected behavior of different users in sharing the road or navigating the intersection.

⊗ LLaVA-RLHF

The objects and features described in the image play a crucial role in ensuring that city traffic flows smoothly and safely. Here are some of the key aspects they contribute to:

1. Traffic lights: The presence of traffic lights at the intersection helps regulate the flow of traffic by controlling the right of way for vehicles and pedestrians. They ensure that traffic moves in an orderly manner, reducing the risk of accidents and congestion. By following the traffic light signals, drivers can safely navigate the intersection and avoid collisions with other vehicles or pedestrians.

2. **Parking meters:** **Parking meters** play a vital role in managing on-street parking in busy urban areas. They help maintain a steady flow of traffic by preventing vehicles from occupying parking spaces for extended periods, which can lead to congestion and reduced accessibility for other drivers. **Parking meters** also generate revenue for the city, which can be used to fund public services and infrastructure projects.

3. Signage: The presence of various signs on the street, such as **no parking signs**, helps communicate important information to drivers and pedestrians. **These signs provide guidance on parking restrictions**, speed limits, and other traffic rules, ensuring that everyone on the road is aware of their responsibilities and the potential consequences of violating traffic regulations.

4. Bike lanes: The presence of bike lanes in the image indicates that the city is promoting alternative modes of transportation, such as cycling. Bike lanes provide a dedicated space for cyclists, reducing the risk of accidents between bikes and cars. They also encourage more people to choose cycling as a means of transportation, which can help reduce traffic congestion and improve air quality in the city.

5. Pedestrian crossings: The presence of pedestrian crossings in the image highlights the importance of ensuring the safety of pedestrians in busy urban areas. These crossings provide designated areas for pedestrians to cross the street safely, reducing the risk of accidents and promoting a more walkable city environment.

In summary, the objects and features described in the image, such as traffic lights, **parking meters**, signage, bike lanes, and pedestrian crossings, play a significant role in ensuring that city traffic flows smoothly and safely. They help regulate traffic flow, manage parking, communicate important information to drivers and pedestrians, promote alternative modes of transportation, and prioritize the safety of pedestrians. By implementing these measures, cities can create a more efficient and sustainable urban environment for all residents.

Figure 7. Qualitative results of different models. **Unreasonable reasoning and extensions** and **scene related hallucinations** are highlighted in color respectively.



Long-form QA



What might be the significance or role of the objects you described in ensuring city traffic flows smoothly?

InstructBLIP

In the image, there are several objects that play a significant role in ensuring city traffic flows smoothly. These objects include streetlights, traffic lights, **cars, bicycles, and pedestrians**. Streetlights provide illumination for drivers and pedestrians during nighttime or low-light conditions, improving visibility and safety on the road. Traffic lights regulate the flow of traffic at intersections, ensuring that vehicles, bicycles, and pedestrians follow the rules of the road and avoid collisions. **Cars and bicycles are essential components of the transportation system, contributing to the movement of people and goods within the city. Pedestrians also play a crucial role in ensuring the smooth flow of traffic by crossing streets safely and adhering to traffic rules.** Overall, these objects work together to create a safe and efficient transportation system in the city.

GPT4-V

The objects in the image are various traffic control devices used in urban settings to manage and regulate the flow of traffic, ensuring safety and order on the roads. Here's their significance and role:

1. Traffic Lights: They control the flow of traffic at intersections, directing vehicles when to stop and go. This prevents collisions and ensures an orderly movement of traffic.
2. Street Signs: Various signs in the image, like "No Right Turn", inform drivers of local traffic laws and regulations they must follow, such as which turns are prohibited at an intersection.
3. Bicycle Lane Markings: The green painted bike lane signifies a dedicated space for cyclists, which helps in reducing conflicts between motor vehicles and bicycles.
4. Pedestrian Signals: Pedestrian traffic lights, although not clearly visible in this image, usually accompany traffic lights to indicate when it is safe for pedestrians to cross the street.
5. Road Markings: The various lines and symbols painted on the road surface guide drivers on how to use the lanes properly, where to stop, and the direction of permissible traffic movement.
6. Parking Signs: These signs inform drivers of the areas where parking is either allowed or prohibited during certain times, which contributes to the reduction of roadside congestion.

Each of these devices plays a crucial role in maintaining a systematic and smooth flow of different modes of transportation within the city streets, minimizing accidents, and ensuring the safety of all road users.

Figure 8. Qualitative results of different models (continued figure). **Unreasonable reasoning and extensions** are highlighted in color.



Long-form QA



Provide a thorough description of the given picture.

⚙️ RLHF-V (This work)

The image depicts a small kitchen with **white walls** and wood cabinets. There is a center island **in the middle of** the kitchen, providing additional counter space. The kitchen is equipped with standard appliances, such as a stove, located on the right side, and a sink on the left.

⚙️ GPT4-V

Certainly! The image showcases an interior view of a kitchen. Here are the details:

Cabinetry: The kitchen is equipped with a combination of **white** overhead and **base cabinets, some of which have a horizontal black line running along their middle**. The design appears modern, with a sleek finish.

Countertop: The countertop appears to be made of a light-colored, possibly marble or granite material. The main countertop extends along one side of the kitchen, while an additional rounded extension provides a potential breakfast or serving area.

Flooring: The floor is checkered with white tiles, and numerous black dots are scattered across each tile, creating a distinct pattern.

Appliances and Fixtures: On the right, there's a silver, freestanding gas range and oven. **Just above it is an exhaust or range hood**. A stainless-steel sink is set into the countertop on the left side of the image, **beneath which there's a dish drying rack attached to the wall**.

Lighting: The ceiling is equipped with multiple **recessed lights**, ensuring the kitchen is well-lit.

Other Details: A **green trash bin** is placed beside the gas range. Various other minor items, like **a bottle and some containers**, can be seen on the countertop. There's also a metallic vertical pole, possibly **a utility or storage pole, standing between the range and the rounded counter extension**.

Walls and Door: The kitchen walls are painted in a muted beige or light brown color. To the far right, there's a wooden door, possibly leading to a pantry or another room.

Ambiance: The kitchen exudes a clean, functional, and organized atmosphere, with a blend of neutral tones and contemporary design elements.

Overall, the kitchen appears to be well-maintained and offers a mix of practicality and style.

Figure 9. Qualitative results of different models. **Hallucinations** and *scene related hallucinations* are highlighted in color respectively.



Long-form QA



Please provide a detailed description of the image. Describe the visual elements, colors, shapes, textures, and any objects or people present along with the overall mood or atmosphere portrayed in the image.



RLHF-V (This work)

The living room is clean and empty of people, with the sun shining in through the **windows**. It features a **white rug** and black furniture that contrasts with the light from the **windows**. There is a doorway leading to a balcony that showcases the railing. The living room contains a couch, **a table**, and two chairs.



GPT4-V

The image showcases an elegantly designed living room.

Starting with the overall atmosphere, the room exudes a warm and inviting ambiance. Natural sunlight floods in from large **windows on the left**, emphasizing the spaciousness of the room and creating a luminous glow.

As for the visual elements:

Colors: The room primarily features neutral tones—creams, browns, and blacks. These muted colors contribute to the room's calming and sophisticated feel. The wooden floor has a polished honey-brown hue which complements the lighter walls.

Furniture: A variety of seating options are available. There are two dark leather chairs **on the left side of the room, facing the window**. In the center, a black leather couch with decorative cushions sits opposite a glass coffee table with metallic legs. **On the far right, there's another seating area with two plush chairs flanking a fireplace.**

Decorative Elements: Above the fireplace is a white mantel adorned with decorative items. On the left wall, a sleek table holds decorative items like **a lamp and a vase** with flowers. **Above this table, a large mirror is hung**, reflecting the room's light and enhancing its spaciousness. **On the adjacent wall**, three framed artworks are neatly aligned, adding an artistic touch to the room.

Textures: The room presents a blend of textures. The plushness of the rug contrasts with the smooth leather of the furniture, the glass surface of the coffee table, and the sleek wooden floor.

Shapes: There are various shapes at play. The rectangular and square forms of the artworks, the **mirror**, and the **windows** are balanced by the **rounded shapes of the seating and the coffee table**.

Plants: A large potted plant sits by the window, adding a touch of nature and a refreshing green to the neutral palette.

The room is a blend of modern and classic elements, providing a sense of comfort while maintaining a refined aesthetic.

Figure 10. Qualitative results of different models. **Hallucinations** and *scene related hallucinations* are highlighted in color respectively.