

MITIGATING OBJECT HALLUCINATION IN LARGE VISION LANGUAGE MODEL WITH HUMAN-FREE REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Vision-Language Models (LVLMs) have excelled in joint visual and language understanding, particularly in generating detailed image captions. However, they still struggle with object hallucination, where non-existent objects are described, especially in long captions. While fine-tuning through supervised learning with enhanced datasets or reinforcement learning from human feedback can alleviate this issue, these methods demand considerable human effort, limiting scalability. This paper addresses this challenge by introducing a human-free framework to mitigate object hallucination in LVLMs for image captioning, utilizing reinforcement learning driven exclusively by automatic natural language processing metrics. We demonstrate that the following framework can effectively mitigate hallucination: (1) caption generation is formulated as a Markov Decision Process (MDP); (2) minimizing hallucination while maintaining caption quality is guided by a reward function, combining a proposed *FIScore* with a penalty on Kullback–Leibler divergence from the pre-trained model; (3) fine-tuning the LVLM within the MDP framework can be performed directly by Proximal Policy Optimization (PPO) with careful attention to architectural details. Extensive experiments demonstrate a significant reduction in hallucination by up to 41% while preserving the caption quality compared to the baseline model, InstructBLIP, on the COCO dataset. This improvement is reflected in consistent gains in object coverage and accuracy across various models and datasets. Notably, our method achieves comparable or superior performance to alternative approaches, all without requiring any human involvement.

1 INTRODUCTION

Large Vision-Language Models (LVLMs) have become increasingly prominent due to their ability to perform joint visual and language understanding tasks Achiam et al. (2023); Alayrac et al. (2022). Among these, image captioning has emerged as a key application where LVLMs consistently outperform smaller models by generating highly detailed and contextually rich captions Dai et al. (2023); Zhu et al. (2023); Li et al. (2023a). Despite these advancements, LVLMs still struggle with a critical challenge: object hallucination Rohrbach et al. (2018b); Biten et al. (2022). This occurs when captions include references to objects that do not exist in the corresponding image, particularly in longer, more detailed descriptions; as shown in Fig. 1. Object hallucination not only undermines the credibility of these models but also hinders their broader application in fields that require high precision, such as autonomous systems and medical imaging.

Addressing object hallucination has been a major focus in recent research efforts Zhou et al. (2023); Li et al. (2023d); Dai et al. (2022); Liu et al. (2024). Early efforts aimed at mitigating this issue in small-scale multimodal pre-trained models focused on reducing object co-occurrence patterns through data augmentation Biten et al. (2022); Rohrbach et al. (2018b); Kim et al. (2023). However, such approaches were considered ineffective for LVLMs Zhou et al. (2023). More recent studies have explored improving dataset quality and applying fine-tuning to LVLMs Gunjal et al. (2023); Li et al. (2023c); Liu et al. (2023a), or using Reinforcement Learning from Human Feedback (RLHF) Sun et al. (2023) to reduce object hallucination. Despite their potential, these methods still face significant challenges, as gathering large volumes of high-quality examples Gunjal et al. (2024); You

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107



(a) Sentence hallucination ratio measured in CHAIR_s. (b) A detailed caption example from COCO dataset for baseline (Base) vs fine-tuned (Our). Bold objects are hallucinated ones by LVLMs.

Figure 1: Quantitative and qualitative comparison between Base (InstructBLIP) and Our: Chart (a) shows a significant 41% reduction in object hallucination on the COCO dataset using Our. Fig. (b) presents an example where the Base model produces a caption with substantial object hallucination, while the Our model provides an accurate description without hallucinated objects.

et al. (2023); Zhang et al. (2024) or obtaining accurate human feedback for RLHF fine-tuning Stienon et al. (2020) remains a time-consuming and labor-intensive process that requires considerable human expertise and effort.

To address these limitations, we propose a human-free framework to mitigate object hallucination in LVLMs for image captioning. Our approach leverages reinforcement learning, guided exclusively by automatic natural language processing (NLP) metrics, eliminating the need for human intervention. The key features of our framework are as follows:

- **Caption Generation as an MDP:** To streamline previous methods and minimize human intervention, we formulate the caption generation task as a Markov Decision Process (MDP), with a reward function incorporating specific automatic NLP metrics to reduce hallucination. By framing image captioning as a reinforcement learning problem, we can effectively address the inherent non-differentiability challenge of optimizing automatic metrics, which are difficult to optimize directly through traditional supervised learning methods.
- **Dedicated Reward Function:** To guide the output generation behavior, we incorporate automatic NLP metrics into the reward function. For hallucination reduction, instead of using the straightforward CHAIR metric Rohrbach et al. (2018b), we introduce FIScore, which provides a better balance between reducing object hallucination and improving object coverage. Additionally, we introduce a Kullback–Leibler (KL) divergence penalty to prevent the policy from diverging too far from the pre-trained model, preserving caption quality without the need for labeled data. Moreover, since metrics like FIScore are computed only at the end of caption generation, which results in sparse rewards, the KL penalty helps densify feedback, making RL optimization more effective. Optionally, when labeled data is available, the reward can easily adopt other quality metrics such as Meteor Banerjee & Lavie (2005) and BERTScore Zhang et al. (2019) to further improve caption quality.
- **Efficient Fine-tuning with PPO:** The proposed framework can be directly optimized using Proximal Policy Optimization (PPO), a popular RL method, to fine-tune the Large Vision-Language Model (LVLM). However, training LVLMs typically requires significant memory and computational resources. To mitigate this, we introduce a compact version of PPO where the policy, value function, and reference model share the same frozen language model, adding only minimal additional training parameters through adapters. These adapters are compatible with recent state-of-the-art fine-tuning techniques for LVLMs such as prompt tuning, ensuring resource-efficient training.

Through extensive experiments, we demonstrate that our method reduces hallucination by up to 41% compared to the baseline model, InstructBLIP, while also improving object coverage and caption quality on the COCO dataset. Additionally, our framework can be easily extended to handle more complex datasets (e.g. Visual Genome) and incorporate existing NLP metrics effectively. Notably,

our approach achieves comparable or superior performance to existing methods, all without relying on human feedback, making it a scalable and efficient solution for enhancing LVLMs in image captioning tasks.

2 RELATED WORK

Large Vision Language Model: The rapid advancements in Large Language Models (LLMs) Touvron et al. (2023); Chung et al. (2022); Touvron et al. (2023) combined with a surge in open-source initiatives, has paved the way for the emergence of extensive vision-language models Liu et al. (2023c); Zhu et al. (2023); Sun et al. (2023); Ye et al. (2023); Bai et al. (2023); Peng et al. (2023). LVLMs seamlessly combine a LLM and a pre-trained visual encoder to form an end-to-end model, aiming to produce contextually relevant text from visual stimuli Zhang et al. (2023a). There are various approaches to effectively achieve this. LLaVA Liu et al. (2023b) introduced the concept of integrating a simple projector during LLM fine-tuning. Chatspot Zhao et al. (2023) follow LLaVA’s model structure, embeds the region of interest into instruction data. GPT4RoI Yu et al. (2023) and Shikra Chen et al. (2023) add grounding tasks to LLaVA structure models and achieve great performance on various tasks. Concurrently, Multimodal-GPT Gong et al. (2023) aims to improve OpenFlamingo’s Alayrac et al. (2022) directive adherence. mPLUG-Owl Ye et al. (2023) suggests a two-step method: first train vision models, and then refine the language model using techniques like LoRA Hu et al. (2021). BLIP2 Li et al. (2023b) and InstructBLIP Dai et al. (2023) presented Q-former-based LVLMs without fine-tuning the LLM but achieving state-of-the-art performance. Our work fine-tunes the InstructBLIP to reduce object hallucination within LVLMs.

Object Hallucination in Vision Language Models: Object hallucination refers to generated descriptions containing objects which are not present in the visual modality Rohrbach et al. (2018b). In small-scale vision language models (VLM), mitigation techniques include fine-grained contrastive learning Zeng et al. (2021) or data augmentation to eliminate co-occurrence patterns Kim et al. (2023). However, training paradigms differ between conventional VLMs and LVLMs. The autoregressive training paradigm in LVLMs poses challenges in implementing VLM hallucination mitigation methods directly Zhang et al. (2023b). Notably, object hallucination is more pronounced and widespread in the long-form descriptions produced by LVLMs compared to the shorter descriptions generated by VLMs. Ongoing research has started to tackle object hallucination in LVLMs, encompassing evaluation and detection approaches Petryk et al. (2024); Li et al. (2023d); Liu et al. (2023a); Dai et al. (2022); Jing et al. (2023); Liu et al. (2023a); Sun et al. (2023), **the development of benchmarks Ben-Kish et al. (2024); Wang et al. (2023)**, hallucination elimination through the construction of higher-quality datasets Gunjal et al. (2023); Li et al. (2023c); You et al. (2023), and the use of supervised learning for fine-tuning Zhou et al. (2023); Zhai et al. (2023) or employ Reinforcement Learning training from Human Feedback (RLHF) Sun et al. (2023); Stiennon et al. (2020) to align different modalities. However, these methods often demand substantial time and labor, particularly in acquiring a large number of high-quality examples. Instead, grounded in reinforcement learning (RL) and automatic metrics, we propose a novel approach. This conceptually distinct method demonstrates efficacy in reducing hallucination and is compatible with various LVLMs, offering a more efficient solution without relying on human effort.

Reinforcement Learning for NLP: Reinforcement Learning (RL) has emerged as a prevalent technique for enhancing language models in a wide range of Natural Language Processing (NLP) tasks, encompassing dialogue Li et al. (2016); Zhou et al. (2017); Jaques et al. (2019); Yi et al. (2019); Jaques et al. (2020), machine translation Wu et al. (2016); Nguyen et al. (2017); Kiegeland & Kreutzer (2021); Bahdanau et al. (2016); Ranzato et al. (2015); Kreutzer et al. (2018), image captioning Rennie et al. (2017); Ren et al. (2017), summarization Stiennon et al. (2020); Paulus et al. (2017); Wu & Hu (2018); Bohm et al. (2019); Ziegler et al. (2019), and text-games Narasimhan et al. (2015); Hausknecht et al. (2020). In this training paradigm, NLP models are optimized through an RL algorithm, wherein the reward signal is derived from either human feedback Kreutzer et al. (2018); Jaques et al. (2020); Stiennon et al. (2020); Ziegler et al. (2019) or NLP evaluation metrics, such as ROUGE for summarization Paulus et al. (2017); Wu & Hu (2018) or BLUE for translation Wu et al. (2016); Nguyen et al. (2017); Kiegeland & Kreutzer (2021). These reward mechanisms enable the models to iteratively improve and fine-tune their performance based on the quality of generated outputs. While RL has proven effective in NLP, its exploration in Vision Large Language Models (LVLMs) for captioning is not well-established. Our work pushes the boundaries in this

direction by leveraging RL to address the challenge of object hallucination in LVLMs. We tackle intricate issues specific to this context, including high computational costs, sparse rewards, and extended temporal horizons.

Finetuning LVLMs with Adapters: Fine-tuning the entire model for Large Vision Language Mode demands extensive memory and computational resources. To address this challenge, various Parameter Efficient Fine-Tuning (PEFT) methods have emerged as cost-effective alternatives. These methods include prompt tuning Lester et al. (2021); Li & Liang (2021); Qin & Eisner (2021), tuning the embedding layer inputs An et al. (2022), tuning hidden states (IA3) Liu et al. (2022), employing Low-rank Adapters (LoRA) Hu et al. (2021); Dettmers et al. (2023), incorporating full layers Housby et al. (2019), tuning biases Zaken et al. (2021), learning weight masks based on Fisher information Sung et al. (2021), and leveraging combinations of these approaches Karimi Mahabadi et al. (2021). In our study, we demonstrate the effectiveness of prompt tuning in addressing the task at hand, while future work will investigate trade-offs with other PEFT methods to further enhance performance.

3 METHODOLOGY

In this session, we will sequentially cover the following topics: (1) casting the caption generation task within the framework of a Markov Decision Process (MDP); (2) defining the dedicated reward function with appropriate automatic metrics; (3) modeling RL networks; (4) fine-tuning the model by solving the MDP through Proximal Policy Optimization (PPO).

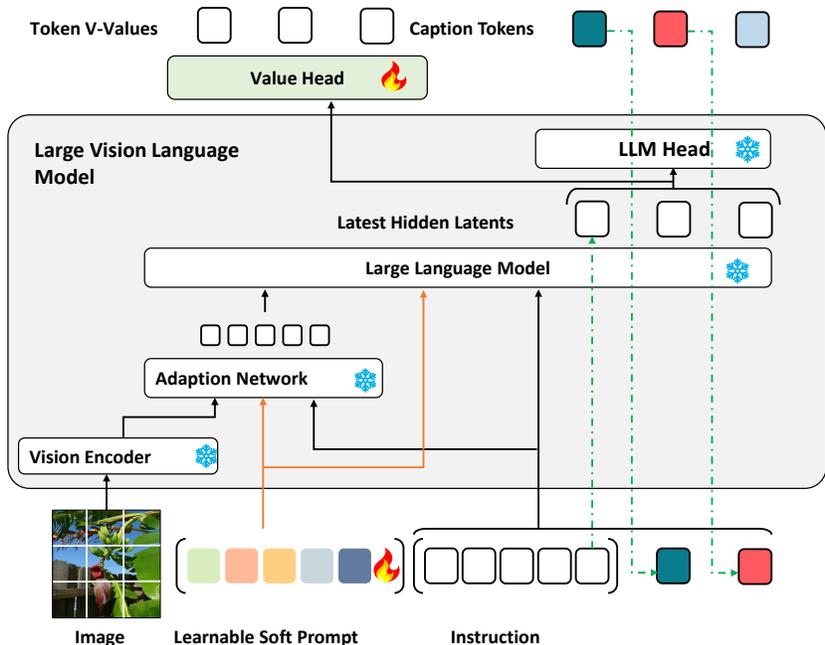


Figure 2: Detailed architecture of our framework. Specifically, the Policy Network is crafted by augmenting the shared LVLm with delicately learnable soft prompts. Meanwhile, the Value Network is formed by replacing the LLM Head with a Linear Value head. Notably, all parameters of the LVLm remain shared and frozen, with only a very small fraction (less than 0.01% LVLm weight) of trainable parameters added to the LVLm for the meticulous modeling of the policy network and value network.

3.1 MARKOV DECISION PROCESS (MDP) FOR IMAGE CAPTIONING

The image captioning task can be effectively framed as an MDP due to its inherent sequential nature, where each token generation is a decision based on the current state. This allows us to utilize RL techniques to optimize caption quality holistically, addressing both local and global aspects

of the generated text. Mathematically, we formulate the image captioning as an MDP denoted by $\langle \mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma, H \rangle$. Each episode in this MDP begins by sampling a datapoint (X, Z, Y) from our dataset $\mathcal{D} = \{(X_i, Z_i, Y_i)\}_{i=1}^N$, where $X \in \mathcal{X}$ represents the text input for LVLMs, $Z \in \mathcal{Z}$ represents the image, and $Y \in \mathcal{Y}$ is the ground truth caption, which can be set to *none* if no ground truth caption is available. The initial state $S_0 = (Z, x_0, \dots, x_m)$ consists the image Z and the text input $X = (x_0, \dots, x_m)$, where $S_0 \in \mathcal{S}$ and the state space $\mathcal{S} = \mathcal{Z} \cup \mathcal{X}$ is defined as the concatenation of images and text inputs. At each time step t , an action $a_t \in \mathcal{A}$, which corresponds to a token from our vocabulary \mathcal{V} , is taken in the environment from a policy (e.g. an LVLm). The transition function $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ deterministically appends an action a_t to the end of the state $S_{t-1} = (Z, x_0, \dots, x_m, a_0, \dots, a_{t-1})$ to form the state S_t . This process continues until the end of the episode $t \leq T \leq H$, either when the current time step t exceeds the horizon H or when an end-of-sentence (EOS) token is generated, resulting in a final state $S_T = (Z, x_0, \dots, x_m, a_0, \dots, a_T)$. At every step, a reward $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{Y} \rightarrow \mathbb{R}^1$ is emitted. This reward may be derived from automated metrics (e.g., CHAIR). Our objective is to maximize the cumulative return represented by the equation:

$$\max_{A=\{a_0 \dots a_T\} \in \mathcal{V}^T} \sum_t \gamma^t \mathcal{R}(S_t, a_t, Y). \quad (1)$$

where γ denotes the discount factor (e.g., 0.99) and A is the generated caption from the LVLm.

3.2 REWARD FUNCTION

To tackle hallucination, the first approach people usually think of is to incorporate $CHAIR_i$ and $CHAIR_s$ Rohrbach et al. (2018a) directly into the reward function. Although $CHAIR$ metrics primarily evaluate precision, they cause models to prioritize precision at the expense of recall. To address this issue, we propose utilizing the $FIScore$. $FIScore$ offers a balanced measure of precision and recall, ensuring that the reward function encourages comprehensive object coverage while maintaining accuracy:

$$FIScore = \frac{2 * Precision * Recall}{Precision + Recall} \quad (2)$$

where $Precision$ is the ratio of correct objects to all predicted objects, and $Recall$ is the ratio of correct objects to all objects in the ground truth. The ground truth objects can be either extracted using an off-the-shelf object detection model (e.g., YOLOv8 Varghese & Sambath (2024)) or obtained directly from the dataset. Predicted objects can be easily extracted from the caption using a method similar to $CHAIR$ Rohrbach et al. (2018a).

The resulting reward function is:

$$\mathcal{R}(S^t, a^t, Y) = \begin{cases} FIScore(S^T, a^T, Y) & \text{if } t = T \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

Optionally, in the setting where ground truth captions are available, two additional metrics can be integrated into the reward function to further enhance caption quality: $Meteor$ Banerjee & Lavie (2005) and $BERTScore$ Zhang et al. (2019). $Meteor$ evaluates the similarity between generated and reference texts (a.k.a ground truth captions) based on n-grams and word order, ensuring structural and lexical alignment. Meanwhile, $BERTScore$ assesses semantic similarity using pre-trained BERT embeddings, capturing underlying meaning accurately. Together, $Meteor$ and $BERTScore$ offer a comprehensive evaluation of caption quality, considering both surface-level and semantic aspects, thereby improving caption relevance to the ground truth.

The enhanced reward function is defined as:

$$\mathcal{R}(S^t, a^t, Y) = \begin{cases} FIScore(S^T, a^T, Y) + \alpha Meteor(S^T, a^T, Y) + \beta BERTScore(S^T, a^T, Y) & \text{if } t = T \\ 0 & \text{otherwise.} \end{cases} \quad (4)$$

Note that the balancing weight α for $Meteor$ should be relatively smaller compared to $FIScore$ and $BERTScore$, as it may encourage shorter captions, especially in datasets with shortened ground truth references.

3.3 MODELING REINFORCEMENT LEARNING NETWORKS

270 To fine-tune the LVLM within our MDP frame-
 271 work, we utilize a policy network, a value net-
 272 work, and a reference network. While the pol-
 273 icy and value networks are essential compo-
 274 nents of RL, the reference network serves as
 275 a proposed teacher network. Its role is to pre-
 276 vent the policy network from deviating too far
 277 from the baseline during training, which is par-
 278 ticularly important for preserving the caption
 279 meaning in the absence of ground truth cap-
 280 tions. Given the intensive computational dem-
 281 ands of conventional LVLM fine-tuning, we
 282 design lightweight and efficient networks. Fig.
 283 3 displays the simplified overview of the frame-
 284 work’s network components. Specifically, each
 285 network builds upon the same frozen LVLM
 286 foundation. The reference network mirrors the
 287 LVLM identically, while the value and policy
 288 networks incorporate slender adapters alongside
 the LVLM. This approach optimizes computa-
 tional resources and is compatible with vari-
 ous state-of-the-art Parameter Efficient Fine-
 Tuning (PEFT) methods Mangrulkar et al. (2022), which rely on adapters.

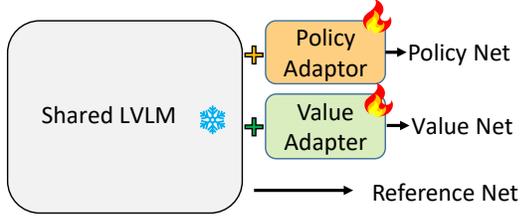


Figure 3: Simplified overview of the framework’s network components: policy network, value network, and reference network. All networks share the same frozen LVLM as its foundation. The reference network mirrors the LVLM identically, while the value and policy networks incorporate a lightweight adapter into the shared LVLM.

289 In this paper, we utilize *Prompt Tuning* to assess the framework’s effectiveness. Prompt Tuning
 290 offers an efficient and flexible method for controlling LVLM behavior. By allowing the model
 291 to remain frozen while refining prompts, this approach reduces computational costs and provides
 292 task-specific adaptability without compromising the model’s generalization capabilities. Specifi-
 293 cally, the LVLM generates captions based on images and instructions in an autoregressive man-
 294 ner. By prefixing a controllable prompt to the instruction, we can influence the model’s behavior
 295 Lester et al. (2021). Mathematically, we adopt a conditional generation perspective, where A
 296 represents a sequence of tokens forming a caption. The captioning process by LVLM is expressed as
 297 $P_\theta(A|X, Z)$, with θ denoting the LVLM’s weight. Prompting enhances the model’s generation of
 298 A by providing additional context, which is achieved by prefixing a token sequence G to the in-
 299 put X . This aids the model in improving the likelihood of generating the ground truth caption Y :
 300 $P_\theta(Y|[G; X], Z)$. Throughout, the model parameters θ remain unchanged. Optimal G selection can
 301 be achieved via manual exploration (*Hard Prompting*) or by representing G with dedicated param-
 302 eters ϕ , refined through gradient descent (*Soft Prompting*). This updates the conditional generation as
 303 $P_{\theta;\phi}(A|[G; X], Z)$, trainable by maximizing reward through backpropagation, with gradient updates
 solely applied to ϕ , i.e., learnable soft prompt.

304 Fig. 2 illustrates the detailed architecture of the Augmented LVLM in our implementation. The Pol-
 305 icy Network $\pi_{\theta;\phi}(A|G, S)$, identical to $P_{\theta;\phi}(A|[G; X], Z)$, is constructed by enhancing the shared
 306 large language model with the delicately *learnable soft prompt*. Concurrently, the Value Network
 307 $V_{\theta;\omega}(S)$ is created by substituting the LLM Head with a Value Head, featuring a single output neu-
 308 ron. The reference network $\pi_\theta^r(A|S)$ remains identical to the original LVLM. Notably, all parameters
 309 of the Large Vision Language Model persist as shared and frozen. Only an extremely small fraction
 310 (approximately 0.01% LVLM weight) of trainable parameters is introduced to meticulously model
 311 the policy network and value network.

313 3.4 FINE-TUNING MODEL BY SOLVING THE MDP

314 Given the MDP and the RL networks, we fine-tune the augmented LVLM, i.e., the policy, using
 315 the on-policy Proximal Policy Optimization (PPO) algorithm Schulman et al. (2017). Formally, this
 316 algorithm trains the policy $\pi_{\theta;\phi}(A|G, S)$ to maximize long-term discounted rewards over generated
 317 captions:

$$318 \mathbb{E}_\pi \left[\sum_{t=0}^T \gamma^t \mathcal{R}(S_t, a_t, Y) \right]. \quad (6)$$

321 We define our V-value and Q-value functions as follows:

$$322 V^\pi(S_t) = \mathbb{E}_{a_t \sim \pi, Y \sim \mathcal{D}} \left[\sum_{\tau=t}^T \gamma^\tau R(S_\tau, a_\tau, Y) \right] \quad (7)$$

$$Q^\pi(S_t, a_t) = \mathbb{E}_{Y \sim \mathcal{D}} R(S_t, a_t, Y) + \gamma \mathbb{E}_{S_{t+1} \sim P} [V^\pi(S_{t+1})]. \quad (8)$$

This leads to the definition of our advantage function:

$$A^\pi(S_t, a_t) = Q^\pi(S_t, a_t) - V^\pi(S_t). \quad (9)$$

We use the previously mentioned value network $V_{\theta;\omega}$ to model the value function, and the mentioned Reference Network $\pi_\theta^r(A|S)$ to generate the initial caption. Following the components defined, we employ the PPO algorithm detailed in Schulman et al. (2017) to fine-tune the policy. To enhance training stability, we approximate the advantage using Generalized Advantage Estimation as outlined in Schulman et al. (2015).

Given a data point tuple (X, Z, Y) and generated caption A from our policy, as the aforementioned environment reward is sequence-level and sparse, we further regularize the reward function using a token-level KL penalty. This penalty ensures the model does not deviate significantly from the original caption generated by $\pi_\theta^r(A|S)$, densifying the reward signal and preserving the quality and meaning of the caption in line with the reference model. This regularization is especially crucial when the ground truth caption Y is unavailable. Formally, the regularized reward function is defined as:

$$\hat{R}(S_t, a_t, Y) = R(S_t, a_t, Y) \quad (10)$$

$$- \lambda \text{KL}(\pi_\theta(a_t | G, S_t) \parallel \pi^r(a_t | S_t)). \quad (11)$$

Here, \hat{R} is the regularized KL reward, KL denotes Kullback–Leibler divergence, and the KL coefficient λ is dynamically adapted, following Ziegler et al. (2019).

4 EXPERIMENT

4.1 EXPERIMENTAL SETTINGS

Datasets: We train and evaluate our method using the COCO dataset, as described by Lin et al. (2014). This dataset serves as a comprehensive collection widely used in tasks such as image recognition, segmentation, and captioning. It encompasses over 300,000 images, covering more than 80 object categories, and is meticulously annotated. For our captioning task, we utilize the Karpathy split Karpathy & Fei-Fei (2015), dividing the dataset into training, validation, and test sets with 82,000, 5,000, and 5,000 images, respectively. Additionally, to prepare the dataset for LVLM fine-tuning, we randomly augment each image with detailed caption instructions. A complete list of instructions is provided in Appendix G.

Implementation detail: We employ InstructBLIP Dai et al. (2023) as our baseline LVLM due to its robust resistance to hallucination compared to others. InstructBLIP adopts the BLIP-2 architecture Li et al. (2023b) and is distinguished by its use of Q-former, a Query Transformer designed for instruction-aware training. In this paper, the vision encoder utilized is ViT-g/14 Fang et al. (2023), while the LLM of choice is Vicuna-7B. During RL fine-tuning, we initialize the model with the pre-trained InstructBLIP checkpoint. Subsequently, we exclusively fine-tune the parameters of our adapters, keeping the image encoder, Q-former, and LLM frozen.

Our experiments are conducted using the Transformers Wolf et al. (2020) and PyTorch Paszke et al. (2019) frameworks. For fine-tuning on the dataset, we employ the same tokenizer as InstructBLIP with vocabulary size \mathcal{V} 32000. Our reward function sets α and β to 0.1 and 1, respectively. The soft prompt length is set to 20. In implementing PPO, we adopt the default parameters of the Stable Baseline API Raffin et al. (2021), with modifications: we gather 4096 transitions and update the PPO loss 5 times for each on-policy step. The γ is set to 0.99. The KL coefficient λ is dynamically adjusted, as described in Ziegler et al. (2019), with a target KL of 0.05. Our batch size is set to 64, and we train the models using the AdamW optimizer with a learning rate of 0.0002, ensuring stable convergence over 50 epochs. We leverage 8 Nvidia A6000 GPUs, employing mixed precision and flash attention mechanisms Dao et al. (2022) to enhance training speed. The fine-tuning process typically requires approximately one day to complete.

4.2 EXPERIMENTAL RESULTS

In this section, we present experimental results that highlight five key points: (1) the occurrence of object hallucination and its amplification in detailed captions; (2) the potential of prompt-tuning

(demonstrated using hard prompting) to mitigate hallucination; (3) the effectiveness of our framework in reducing hallucination while preserving or even enhancing caption quality when ground truth captions are available, compared to baseline models and alternative methods; and (4) the framework’s robustness when applied to more complex datasets.

The occurrence of object hallucination and its amplification in detailed captions: We begin by conducting an experiment aimed at demonstrating the presence of object hallucination and its amplification with detailed captions. We design instructions to generate short and long captions using two baseline models: InstructBLIP and mPLUG-Owl. Tab. 1 illustrates the object hallucination measured by CHAIRs across various caption types with specific input prompts on the COCO test set. The results indicate that LVLMS experience object hallucination for both short and long captions, with the issue being more pronounced for longer captions. Notably, InstructBLIP exhibits less hallucination with short captions; however, the problem amplifies significantly, around ten times, with longer sentences. Both models show similarly high rates of hallucination in long captions demonstrating the severity of the problem.

Type	Prompt	InstructBLIP		mPLUG-Owl	
		$CHAIR_i(\%) \downarrow$	$CHAIR_s(\%) \downarrow$	$CHAIR_i(\%) \downarrow$	$CHAIR_s(\%) \downarrow$
Short	Generate a short caption of the image.	2.43	3.13	22.81	60.55
	Create a textual summary for the image.	4.95	6.51	22.98	61.33
Long	Provide a detailed description of the image.	27.01	60.91	26.03	71.39
	Create a detailed textual summary for the image.	25.80	59.11	24.25	66.31

Table 1: Object Hallucination, gauged by $CHAIR_s$ and $CHAIR_i$ metrics, across diverse caption types paired with specific input prompts in the COCO test set. These prompts are designed to elicit both short and long captions. Two distinct methods are illustrated: InstructBLIP and mPLUG-Owl.

Prompt tuning in mitigating hallucination:

We have meticulously curated a series of hard prompts intended to be incorporated at the beginning of input instructions, aimed at minimizing object hallucination in the model’s generated captions. Each prompt is meticulously designed to address specific sources of object hallucination, strategically guiding the model away from potential pitfalls. The comprehensive list of prompts is provided in Appendix H. During the testing phase, we employ a randomized approach by selecting a single hard prompt to prefix each sample instruction. We conduct captioning using the InstructBLIP baseline model with prefixed instructions. The reported performance metrics reflect the average performance across these instances, focusing particularly on $CHAIR$ evaluations as shown in Tab. 4 under the label *Hard Prompting*.

In comparison to InstructBLIP, we observe that hard prompting can mitigate object hallucination by reducing $CHAIR_i$ and $CHAIR_s$ from 25.8 to 20.9 (-4.9%) and 59.1 to 45.1 (-14%) respectively. This highlights the effectiveness of prompt tuning as a method to reduce object hallucination.

Performance of our framework: Based on the observed effectiveness of hard prompting, we fine-tuned the InstructBLIP model using a learnable soft prompt within our framework to optimize prompt selection. In Table 4, we present the performance of our proposed method compared to various baselines. The first row represents hallucination of the state-of-the-art LVLMS before fine-tuning: mPLUG-Owl Li et al. (2022), LLaVA Liu et al. (2023b), InstructBLIP Li et al. (2023a). We collected several fine-tuning approaches on the baseline InstructBLIP in the second rows as presented by Zhou et al. (2023).

Method	$CHAIR_i(\%) \downarrow$	$CHAIR_s(\%) \downarrow$
mPLUG-Owl	26.2	70.5
LLaVA	22.5	62.7
InstructBLIP	25.8	59.1
Teacher	7.5	36.4
CoT	7.8	35.7
Greedy-Decoding	7.8	35.5
GPT-Ensemble	13.0	51.0
GPT-Teacher	7.8	32.0
Hard Prompting	20.9	45.1
Our	6.8	17.8

Figure 4: Performance of Object Hallucination. The first row showcases non-fine-tuned LVLMS baselines. The second row features fine-tuning methods referenced in Zhou et al. (2023). The third row illustrates our Hard Prompting on baseline InstructBLIP, while the last row demonstrates our Soft Prompt fine-tuning using our RL framework.

The results demonstrate that our proposed method consistently outperforms all non-fine-tuning baselines across hallucination metrics. Remarkably, our approach enhances $CHAIR_i$ by +18.9% and $CHAIR_s$ by +41.3% compared to the baseline InstructBLIP and notably surpasses the performance of *Hard Prompting*. Among fine-tuning approaches, we achieved the top ranking on $CHAIR_s$ and second place on $CHAIR_i$, with a very marginal difference compared to the best-performing model, LURE.

Additionally, our method is able to maintain or enhance caption quality across various metrics. Table 5 presents our results. The row for **Our** demonstrates the use of $FIScore$ and KL divergence, maintaining performance comparable to the base model, InstructBLIP. There is a slight increase in $SPICE$, $BLUE$, and $BERTScore$, which we attribute to the generated captions being more factual, concise, and focused, resulting in shorter and more precise outputs. When ground truth captions are available, incorporating Meteor and $BERTScore$, as in **Our-Enhance**, significantly improves caption quality. It is evident that **Our-Enhance** significantly improves captioning performance across $SPICE$, $BLUE$, and $BERTScore$, surpassing all previous baselines.

Extend evaluations to complex dataset: We conducted additional evaluations using the Visual Genome dataset and the CCEval metric as outlined in Halle-switch Zhai et al. (2023). These evaluations allowed us to explore the model’s performance in more complex scenarios, where captions typically contain a denser array of objects, potentially increasing the likelihood of hallucination. The result is shown in Fig. 6.

Interestingly, the LLaVa13B model, despite being a stronger generative model, shows more hallucinations in both CCEval-i and CCEval-s scores compared to LLaVa7B. Examining the generated captions shows that this is due to LLaVa13B’s tendency to generate more imaginative content, indicating that while increased model capability can enhance creativity, it may also lead to more hallucinations. Therefore, guiding the model to prioritize factual accuracy is essential.

Fig. 6 also clearly shows the effectiveness of our model in reducing object hallucinations, with significantly lower $CCEVal_i$ (object-level) score and achieve the best $CCEVal_s$ (caption-level) score among baseline models. Although the improvement in $CCEVal_s$ is marginal, this is likely due to the higher object density in Visual Genome Images, which increases the risk of hallucination and makes it more challenging to eliminate hallucinations entirely. Nonetheless, our model demonstrates robustness and adaptability in handling complex captioning tasks, confirming its effectiveness beyond the COCO dataset.

4.3 ABLATION STUDY

Effectiveness of $FIScore$: The $FIScore$ plays a crucial role in ensuring the recall of generated captions. Fig. 7 provides a comparison between using the $FIScore$ instead of $CHAIR$ in the reward. It is evident that employing $CHAIR$ directly has a detrimental effect, significantly re-

Method	Captioning Quality		
	$SPICE$ (%) \uparrow	$BLEU$ (%) \uparrow	$BERTScore$ (%) \uparrow
mPLUG-Owl	12.5	2.7	87.40
LLaVA	13.5	3.0	87.83
InstructBLIP	10.9	1.1	85.81
Hard Prompt	11.1	1.0	85.9
Our	11.0	1.5	86.86
Our-Enhance	14.6	6.6	90.42

Figure 5: Captioning quality is evaluated using NLP metrics, comparing our approach to other methods. **Our** uses only $FIScore$ and KL divergence, while **Our-Enhance** incorporates additional metrics: Meteor and $BERTScore$.

Model	$CCEVal_i$ \downarrow	$CCEVal_s$ \downarrow
LLaVA7B	72.00	19.7
LLaVA13B	79.00	23.8
InstructBlip7B	72.00	22.30
Our	27.0	19.6

Figure 6: The performance of our method on the Genome dataset.

Reward	Pre (%) \uparrow	Rec (%) \uparrow	$CHAIR_i$ (%) \downarrow	$CHAIR_s$ (%) \downarrow
Base	72.9	71.3	27.1	60.9
$CHAIR$	93.7	20.6	6.3	14.4
$FIScore$	93.2	70.2	6.8	18.8

Figure 7: Comparison of $Precision$ (Pre) and $Recall$ (Rec) between using $CHAIR$ and $FIScore$ in the reward function.

486 ducing the recall. This outcome can be attributed to the sole emphasis on precision without due
 487 consideration for recall. The *FIScore* addresses this issue by incentivizing the model to maintain
 488 high recall, thus preserving a comprehensive coverage of ground truth objects.

489 Ablation on Incorporating NLP

490 **Metrics:** Fig. 8 illustrates the impact of using different automatic metrics.
 491 The baseline model shows high object hallucination with 25.8%
 492 under *CHAIR_i*. Incorporating the *FIScore* significantly reduces hallucination
 493 down to 6.8% while maintaining comparable *BERTScore* and
 494 *BLEU* score to the baseline. Adding *BERTScore* and *Meteor* to the reward function further
 495 enhances caption quality, achieving 92.42 in *BERTScore* and 6.6 in *BLEU* on the COCO test
 496 dataset. This ablation study highlights the effectiveness of each component, particularly the
 497 *FIScore*'s role in reducing hallucination, and the additional benefits of *BERTScore* and
 498 *Meteor* for optimizing caption quality when reference captions are available.

Base	FIScore	BERTScore	Meteor	<i>CHAIR_i</i>	<i>BERTScore</i>	<i>BLEU</i>
✓				25.8	85.81	1.1
✓	✓			6.8	86.86	1.5
✓	✓	✓		6.9	90.51	1.8
✓	✓	✓	✓	6.9	90.42	6.6

499 Figure 8: The ablation studies examining the impact of *BERTScore* and *Meteor* metrics on the COCO test set
 500 *BERTScore* and *Meteor* to the reward function further enhances caption quality, achieving 92.42 in
 501 *BERTScore* and 6.6 in *BLEU* on the COCO test dataset. This ablation study highlights the effective-
 502 ness of each component, particularly the *FIScore*'s role in reducing hallucination, and the additional
 503 benefits of *BERTScore* and *Meteor* for optimizing caption quality when reference captions are avail-
 504 able.

505 5 DISCUSSION

506 **On the scalability and Computational Resources:** Our framework performs LVLm fine-tuning
 507 by leveraging automatic NLP metrics, significantly reducing the reliance on human effort, thus en-
 508 hancing scalability. The quality of the fine-tuned model depends on automatic metrics like *FIScore*.
 509 As more advanced hallucination metrics are developed, our framework can easily integrate them
 510 without major changes.

511 During development, we recognized the significant GPU demands of fine-tuning LVLms. To ad-
 512 dress this, we designed the framework with efficiency at its core, eliminating network duplication
 513 and leveraging the PEFT approach. It is worth noting that combining mixed precision Micikevicius
 514 et al. (2017) with efficient attention mechanisms (e.g. xformers Lefaudeux et al. (2022)) and ad-
 515 vanced distributed training methods (e.g. Accelerate Gugger et al. (2022)) synergistically supports
 516 our framework's implementation. With adequate GPU resources, our approach is highly suitable.
 517 However, future work could explore prediction-time adaptations, such as prompt engineering, to
 518 scale models even larger and provide more accessibility to hobbyist researchers. Larger LVLms,
 519 with stronger prompt-following capabilities, are especially likely to benefit from these methods.

520 **On the detailed caption length:** Our results demonstrate a significant reduction in hallucinations,
 521 but we observed a minor side effect: the average caption length is shorter than the baseline (85 to-
 522 kens compared to 110 tokens). A closer examination revealed that the model's emphasis on factual
 523 content leads to the omission of imaginative elements, resulting in shorter captions. Our experi-
 524 ments indicate that penalizing shorter captions (in the reward function) can increase their length to
 525 approximately 105 tokens. Unfortunately, this adjustment also raises the hallucination rate to 7.8%.
 526 This suggests a trade-off between caption length and hallucination rates that we should be aware of.
 527 Balancing these factors is crucial for optimizing performance based on specific needs.

528 6 CONCLUSION

529 In conclusion, this paper addresses the persistent challenge of object hallucination in LVLms for
 530 image captioning, especially in detailed descriptions. Traditional fine-tuning methods, while ef-
 531 fective, face scalability issues due to substantial human effort requirements. To overcome this, we
 532 propose a novel framework that leverages reinforcement learning (RL) with automatic natural lan-
 533 guage processing metrics within an MDP framework. This approach minimizes object hallucination
 534 while preserving caption quality, achieved through careful architectural design and a tailored reward
 535 function. Our framework effectively reduces hallucination compared to the baseline model, In-
 536 structBLIP, while maintaining consistent object coverage and caption quality. With its emphasis on
 537 speed and memory efficiency, the framework offers practical scalability and represents a significant
 538 advancement in improving the reliability of LVLms for image captioning.

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545
546 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel
547 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language
548 model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–
23736, 2022.
- 549
550 Shengnan An, Yifei Li, Zeqi Lin, Qian Liu, Bei Chen, Qiang Fu, Weizhu Chen, Nanning Zheng, and
551 Jian-Guang Lou. Input-tuning: Adapting unfamiliar inputs to frozen pretrained models. *arXiv
552 preprint arXiv:2203.03131*, 2022.
- 553
554 Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propo-
555 sitional image caption evaluation. In *Computer Vision–ECCV 2016: 14th European Confer-
556 ence, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part V 14*, pp. 382–398.
Springer, 2016.
- 557
558 Dzmitry Bahdanau, Philemon Brakel, Kelvin Xu, Anirudh Goyal, Ryan Lowe, Joelle Pineau, Aaron
559 Courville, and Yoshua Bengio. An actor-critic algorithm for sequence prediction. *arXiv preprint
560 arXiv:1607.07086*, 2016.
- 561
562 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang
563 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities.
arXiv preprint arXiv:2308.12966, 2023.
- 564
565 Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved
566 correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic
567 evaluation measures for machine translation and/or summarization*, pp. 65–72, 2005.
- 568
569 Assaf Ben-Kish, Moran Yanuka, Morris Alper, Raja Giryes, and Hadar Averbuch-Elor. Mitigating
570 open-vocabulary caption hallucinations. In *Proceedings of the 2024 Conference on Empirical
571 Methods in Natural Language Processing*, pp. 22680–22698, 2024.
- 572
573 Ali Furkan Biten, Lluís Gómez, and Dimosthenis Karatzas. Let there be a clock on the beach:
574 Reducing object hallucination in image captioning. In *Proceedings of the IEEE/CVF Winter
575 Conference on Applications of Computer Vision*, pp. 1381–1390, 2022.
- 576
577 Florian Bohm, Yang Gao, Christian M Meyer, Ori Shapira, Ido Dagan, and Iryna Gurevych. Bet-
578 ter rewards yield better summaries: Learning to summarise without references. *arXiv preprint
579 arXiv:1909.01214*, 2019.
- 580
581 Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing
582 multimodal llm’s referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023.
- 583
584 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,
585 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned lan-
586 guage models. *arXiv preprint arXiv:2210.11416*, 2022.
- 587
588 W Dai, J Li, D Li, AMH Tiong, J Zhao, W Wang, B Li, P Fung, and S Hoi. Instructblip: To-
589 wards general-purpose vision-language models with instruction tuning. arxiv 2023. *arXiv preprint
590 arXiv:2305.06500*, 2023.
- 591
592 Wenliang Dai, Zihan Liu, Ziwei Ji, Dan Su, and Pascale Fung. Plausible may not be faithful: Probing
593 object hallucination in vision-language pre-training. *arXiv preprint arXiv:2210.07688*, 2022.
- 589
590 Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and
591 memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Process-
592 ing Systems*, 2022.
- 593
Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning
of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.

- 594 Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong
595 Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale.
596 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
597 19358–19369, 2023.
- 598
599 Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu,
600 Wenwei Zhang, Ping Luo, and Kai Chen. Multimodal-gpt: A vision and language model for
601 dialogue with humans. *arXiv preprint arXiv:2305.04790*, 2023.
- 602
603 Sylvain Gugger, Lysandre Debut, Thomas Wolf, Philipp Schmid, Zachary Mueller, Sourab Man-
604 grulkar, Marc Sun, and Benjamin Bossan. Accelerate: Training and inference at scale made sim-
605 ple, efficient and adaptable. <https://github.com/huggingface/accelerate>, 2022.
- 606
607 Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and preventing hallucinations in large vision
608 language models. *arXiv preprint arXiv:2308.06394*, 2023.
- 609
610 Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and preventing hallucinations in large vision
611 language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38,
612 pp. 18135–18143, 2024.
- 613
614 Matthew Hausknecht, Prithviraj Ammanabrolu, Marc-Alexandre Côté, and Xingdi Yuan. Interac-
615 tive fiction games: A colossal adventure. In *Proceedings of the AAAI Conference on Artificial
616 Intelligence*, volume 34, pp. 7903–7910, 2020.
- 617
618 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, An-
619 drea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp.
620 In *International Conference on Machine Learning*, pp. 2790–2799. PMLR, 2019.
- 621
622 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
623 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint
624 arXiv:2106.09685*, 2021.
- 625
626 Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza, Noah
627 Jones, Shixiang Gu, and Rosalind Picard. Way off-policy batch deep reinforcement learning of
628 implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456*, 2019.
- 629
630 Natasha Jaques, Judy Hanwen Shen, Asma Ghandeharioun, Craig Ferguson, Agata Lapedriza, Noah
631 Jones, Shixiang Shane Gu, and Rosalind Picard. Human-centric dialog training via offline rein-
632 forcement learning. *arXiv preprint arXiv:2010.05848*, 2020.
- 633
634 Liqiang Jing, Ruosen Li, Yunmo Chen, Mengzhao Jia, and Xinya Du. Faithscore: Evaluating hallu-
635 cinations in large vision-language models. *arXiv preprint arXiv:2311.01477*, 2023.
- 636
637 Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. Compacter: Efficient low-rank
638 hypercomplex adapter layers. *Advances in Neural Information Processing Systems*, 34:1022–
639 1035, 2021.
- 640
641 Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descrip-
642 tions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
643 3128–3137, 2015.
- 644
645 Samuel Kiegeand and Julia Kreutzer. Revisiting the weaknesses of reinforcement learning for
646 neural machine translation. *arXiv preprint arXiv:2106.08942*, 2021.
- 647
648 Jae Myung Kim, A Koepke, Cordelia Schmid, and Zeynep Akata. Exposing and mitigating spurious
649 correlations for cross-modal retrieval. In *Proceedings of the IEEE/CVF Conference on Computer
650 Vision and Pattern Recognition*, pp. 2584–2594, 2023.
- 651
652 Julia Kreutzer, Shahram Khadivi, Evgeny Matusov, and Stefan Riezler. Can neural machine trans-
653 lation be improved with user feedback? *arXiv preprint arXiv:1804.05958*, 2018.

- 648 Benjamin Lefaudeux, Francisco Massa, Diana Liskovich, Wenhan Xiong, Vittorio Caggiano, Sean
649 Naren, Min Xu, Jieru Hu, Marta Tintore, Susan Zhang, Patrick Labatut, Daniel Haziza, Luca
650 Wehrstedt, Jeremy Reizenstein, and Grigory Sizov. xformers: A modular and hackable trans-
651 former modelling library. <https://github.com/facebookresearch/xformers>,
652 2022.
- 653 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt
654 tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- 655
656 Chenliang Li, Haiyang Xu, Junfeng Tian, Wei Wang, Ming Yan, Bin Bi, Jiabo Ye, Hehong Chen,
657 Guohai Xu, Zheng Cao, et al. mplug: Effective and efficient vision-language learning by cross-
658 modal skip-connections. *arXiv preprint arXiv:2205.12005*, 2022.
- 659
660 Jiwei Li, Will Monroe, Alan Ritter, Michel Galley, Jianfeng Gao, and Dan Jurafsky. Deep reinforce-
661 ment learning for dialogue generation. *arXiv preprint arXiv:1606.01541*, 2016.
- 662
663 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-
664 image pre-training with frozen image encoders and large language models. *arXiv preprint*
665 *arXiv:2301.12597*, 2023a.
- 666
667 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-
668 image pre-training with frozen image encoders and large language models. *arXiv preprint*
669 *arXiv:2301.12597*, 2023b.
- 670
671 Lei Li, Yuwei Yin, Shicheng Li, Liang Chen, Peiyi Wang, Shuhuai Ren, Mukai Li, Yazheng Yang,
672 Jingjing Xu, Xu Sun, et al. M₃ it: A large-scale dataset towards multi-modal multilingual in-
673 struction tuning. *arXiv preprint arXiv:2306.04387*, 2023c.
- 674
675 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv*
676 *preprint arXiv:2101.00190*, 2021.
- 677
678 Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating
679 object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023d.
- 680
681 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
682 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer*
683 *Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014,*
684 *Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- 685
686 Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large
687 multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023a.
- 688
689 Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiutian Zhao, Ke Wang, Liping Hou,
690 Rongjun Li, and Wei Peng. A survey on hallucination in large vision-language models. *arXiv*
691 *preprint arXiv:2402.00253*, 2024.
- 692
693 Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and
694 Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context
695 learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022.
- 696
697 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023b.
- 698
699 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv*
700 *preprint arXiv:2304.08485*, 2023c.
- 701
702 Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin
703 Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>, 2022.
- 704
705 Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia,
706 Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision
707 training. *arXiv preprint arXiv:1710.03740*, 2017.

- 702 Karthik Narasimhan, Tejas Kulkarni, and Regina Barzilay. Language understanding for text-based
703 games using deep reinforcement learning. *arXiv preprint arXiv:1506.08941*, 2015.
704
- 705 Khanh Nguyen, Hal Daumé III, and Jordan Boyd-Graber. Reinforcement learning for bandit neural
706 machine translation with simulated human feedback. *arXiv preprint arXiv:1707.07402*, 2017.
- 707 Kishore Papineni, Salim Roukos, Todd Ward, and W^{BLEU} Zhu. A method for automatic evalu-
708 ation of machine translation”. *the Proceedings of ACL-2002, ACL, Philadelphia, PA, July 2002*,
709 2001.
- 710 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
711 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-
712 performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
713
- 714 Romain Paulus, Caiming Xiong, and Richard Socher. A deep reinforced model for abstractive
715 summarization. *arXiv preprint arXiv:1705.04304*, 2017.
- 716 Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu
717 Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint*
718 *arXiv:2306.14824*, 2023.
- 719
- 720 Suzanne Petryk, David M Chan, Anish Kachinthaya, Haodi Zou, John Canny, Joseph E Gonzalez,
721 and Trevor Darrell. Aloha: A new measure for hallucination in captioning models. *arXiv preprint*
722 *arXiv:2404.02904*, 2024.
- 723 Guanghui Qin and Jason Eisner. Learning how to ask: Querying lms with mixtures of soft prompts.
724 *arXiv preprint arXiv:2104.06599*, 2021.
- 725
- 726 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
727 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of*
728 *Machine Learning Research*, 22(268):1–8, 2021. URL [http://jmlr.org/papers/v22/](http://jmlr.org/papers/v22/20-1364.html)
729 [20-1364.html](http://jmlr.org/papers/v22/20-1364.html).
- 730 Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. Sequence level train-
731 ing with recurrent neural networks. *arXiv preprint arXiv:1511.06732*, 2015.
- 732
- 733 Zhou Ren, Xiaoyu Wang, Ning Zhang, Xutao Lv, and Li-Jia Li. Deep reinforcement learning-based
734 image captioning with embedding reward. In *Proceedings of the IEEE conference on computer*
735 *vision and pattern recognition*, pp. 290–298, 2017.
- 736
- 737 Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. Self-critical
738 sequence training for image captioning. In *Proceedings of the IEEE conference on computer*
739 *vision and pattern recognition*, pp. 7008–7024, 2017.
- 740
- 741 Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object
742 hallucination in image captioning. In *Proceedings of the 2018 Conference on Empirical Methods*
743 *in Natural Language Processing*, 2018a.
- 744
- 745 Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object
746 hallucination in image captioning. *arXiv preprint arXiv:1809.02156*, 2018b.
- 747
- 748 John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-
749 dimensional continuous control using generalized advantage estimation. *arXiv preprint*
750 *arXiv:1506.02438*, 2015.
- 751
- 752 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
753 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 754
- 755 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,
756 Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances*
757 *in Neural Information Processing Systems*, 33:3008–3021, 2020.
- 758
- 759 Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan,
760 Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with
761 factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023.

- 756 Yi-Lin Sung, Varun Nair, and Colin A Raffel. Training neural networks with fixed sparse masks.
757 *Advances in Neural Information Processing Systems*, 34:24193–24205, 2021.
758
- 759 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
760 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
761 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 762 Rejin Varghese and M Sambath. Yolov8: A novel object detection algorithm with enhanced perfor-
763 mance and robustness. In *2024 International Conference on Advances in Data Engineering and*
764 *Intelligent Computing Systems (ADICS)*, pp. 1–6. IEEE, 2024.
- 765 Junyang Wang, Yiyang Zhou, Guohai Xu, Pengcheng Shi, Chenlin Zhao, Haiyang Xu, Qinghao Ye,
766 Ming Yan, Ji Zhang, Jihua Zhu, et al. Evaluation and analysis of hallucination in large vision-
767 language models. *arXiv preprint arXiv:2308.15126*, 2023.
768
- 769 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
770 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick
771 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gug-
772 ger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art
773 natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in*
774 *Natural Language Processing: System Demonstrations*, pp. 38–45, Online, October 2020. As-
775 sociation for Computational Linguistics. URL [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/2020.emnlp-demos)
776 [2020.emnlp-demos](https://www.aclweb.org/anthology/2020.emnlp-demos). 6.
- 777 Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey,
778 Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google’s neural machine trans-
779 lation system: Bridging the gap between human and machine translation. *arXiv preprint*
780 *arXiv:1609.08144*, 2016.
- 781 Yuxiang Wu and Baotian Hu. Learning to extract coherent summary via deep reinforcement learn-
782 ing. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
783
- 784 Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen
785 Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models
786 with multimodality. *arXiv preprint arXiv:2304.14178*, 2023.
- 787 Sanghyun Yi, Rahul Goel, Chandra Khatri, Alessandra Cervone, Tagyoung Chung, Behnam He-
788 dayatnia, Anu Venkatesh, Raefer Gabriel, and Dilek Hakkani-Tur. Towards coherent and engag-
789 ing spoken dialog response generation using automatic conversation evaluators. *arXiv preprint*
790 *arXiv:1904.13015*, 2019.
- 791 Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao,
792 Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity.
793 *arXiv preprint arXiv:2310.07704*, 2023.
794
- 795 Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang,
796 and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv*
797 *preprint arXiv:2308.02490*, 2023.
- 798 Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning
799 for transformer-based masked language-models. *arXiv preprint arXiv:2106.10199*, 2021.
800
- 801 Yan Zeng, Xinsong Zhang, and Hang Li. Multi-grained vision language pre-training: Aligning texts
802 with visual concepts. *arXiv preprint arXiv:2111.08276*, 2021.
- 803 Bohan Zhai, Shijia Yang, Chenfeng Xu, Sheng Shen, Kurt Keutzer, and Manling Li. Halle-switch:
804 Controlling object hallucination in large vision language models. *arXiv e-prints*, pp. arXiv–2310,
805 2023.
- 806 Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks:
807 A survey. *arXiv preprint arXiv:2304.00685*, 2023a.
808
- 809 Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A Smith. How language model
hallucinations can snowball. *arXiv preprint arXiv:2305.13534*, 2023b.

- 810 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. BertScore: Evaluat-
811 ing text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- 812
- 813 Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li, Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin,
814 Tong Luo, Yaqian Li, Shilong Liu, et al. Recognize anything: A strong image tagging model.
815 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
816 1724–1732, 2024.
- 817 Liang Zhao, En Yu, Zheng Ge, Jinrong Yang, Haoran Wei, Hongyu Zhou, Jianjian Sun, Yuang Peng,
818 Runpei Dong, Chunrui Han, et al. Chatspot: Bootstrapping multimodal llms via precise referring
819 instruction tuning. *arXiv preprint arXiv:2307.09474*, 2023.
- 820
- 821 Li Zhou, Kevin Small, Oleg Rokhlenko, and Charles Elkan. End-to-end offline goal-oriented dialog
822 policy learning via policy gradient. *arXiv preprint arXiv:1712.02838*, 2017.
- 823
- 824 Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit
825 Bansal, and Huaxiu Yao. Analyzing and mitigating object hallucination in large vision-language
826 models. *arXiv preprint arXiv:2310.00754*, 2023.
- 827
- 828 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: En-
829 hancing vision-language understanding with advanced large language models. *arXiv preprint
arXiv:2304.10592*, 2023.
- 830
- 831 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul
832 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv
preprint arXiv:1909.08593*, 2019.
- 833
- 834

835 A OPEN-VOCABULARY BENCHMARK

836

837

838 In our approach, we evaluate both a closed dataset (COCO) and an open-vocabulary dataset (Vi-
839 sual Genome). For COCO, we selected CHAIR due to its tailored design for this dataset, ensur-
840 ing reliable and consistent results. For Visual Genome, we opted for CCEVAL, which builds on
841 CHAIR’s methodology by incorporating large language models (LLMs) to better capture objects in
842 open-vocabulary settings, particularly in the context of the Visual Genome dataset. Notably, open-
843 vocabulary benchmarks can also be leveraged to evaluate the framework in broader applications.
844 Specifically, the study Mitigating Open-Vocabulary Caption Hallucinations introduces the Open-
845 Chair benchmark, an extension of CHAIR that accommodates a broader object vocabulary than
846 COCO. OpenChair proposes an evaluation method using LLMs to identify hallucinated objects,
847 providing complementary insights for experiments beyond the COCO dataset. Similarly, ALOHa
848 highlights CHAIR’s limitations due to its reliance on string matching for a fixed object set. While
849 CHAIR performs well for COCO, its applicability is limited in open-vocabulary contexts. To over-
850 come this, ALOHa employs LLMs to detect objects in more general settings, enhancing its adapt-
851 ability.

852 It is important to note that CCEVAL, OpenChair, and ALOHa all address the limitations of CHAIR
853 and converge on a shared approach: leveraging LLMs to enable more generalized and versatile
854 applications across diverse datasets.

855

856 B MOTIVATION OF USING REINFORCEMENT LEARNING

857

858 Our motivation for employing RL stems from the need to minimize human effort while ensuring
859 effectively reduct hallucination.
860 Traditional approaches to mitigating hallucinations often require identifying specific sources of hal-
861 lucination and designing targeted strategies to counter them. While effective, these methods are
862 labor-intensive. Data-driven alternatives like supervised learning provide some level of automation
863 but rely heavily on labeled datasets, which still require significant human input for data annotation
and curation—an increasingly costly and time-intensive process, particularly for large-scale models.

864 In contrast, reinforcement learning in the literature not only demonstrates strong alignment capabilities
865 for LVLMs in tasks like image captioning but also offers a promising path to automation by
866 significantly reducing the need for explicit labels (e.g., relying only on simple binary feedback for
867 reward modeling). We are motivated to push this approach to its limits by completely eliminating
868 human-labeled data, fully leveraging RL’s potential through the exclusive use of automatic metrics
869 to reduce hallucinations. These metrics are gradually improving in their alignment with human
870 feedback in terms of both accuracy and reliability. Our approach enables the model to iteratively
871 refine its outputs based solely on automatic feedback, providing an efficient and scalable solution
872 that aligns with the trend toward larger LVLMs.

874 C DESCRIPTIONS OF EVALUATION METRICS

876 **BLEU:** BLEU (Bilingual Evaluation Understudy) is a metric employed for assessing the quality of
877 machine-generated translations by comparing them to one or more reference translations. Derived
878 from the concept of precision in n-grams—consecutive sequences of n words—BLEU quantifies the
879 extent to which the generated translation aligns with the reference translations in terms of n-gram
880 overlap Papineni et al. (2001)

882 **BERTScore:** BERTScore is a technique designed to assess the performance of natural language
883 generation or summarization systems, as introduced by Zhang et al. (2019). This method gauges the
884 similarity between a reference text and a generated text by leveraging contextualized embeddings
885 derived from BERT (Bidirectional Encoder Representations from Transformers).

886 **SPICE:** SPICE (Semantic Propositional Image Caption Evaluation) Anderson et al. (2016) is em-
887 ployed to assess the quality of image captions by evaluating both the semantic content and precision
888 of the generated captions in comparison to reference captions. This metric operates on the hypoth-
889 esis that semantic propositional content plays a crucial role in human caption evaluation. SPICE
890 introduces an automated caption evaluation method defined over scene graphs, aiming to capture
891 the intricacies of semantic representation in image captions.

892 **METEOR:** METEOR (Metric for Evaluation of Translation with Explicit ORdering) Banerjee &
893 Lavie (2005) serves as an evaluation metric for machine translation output. This metric calculates the
894 harmonic mean of unigram precision and recall, with recall carrying greater weight than precision.
895 Unlike other metrics, METEOR incorporates additional features such as stemming and synonymy
896 matching, in addition to the standard exact word matching. Its design addresses certain issues iden-
897 tified in the widely used BLEU metric, aiming to improve correlation with human judgment at the
898 sentence or segment level. Notably, METEOR focuses on sentence-level correlation, diverging from
899 BLEU, which seeks correlation at the corpus level.

900 D LARGE VISION-LANGUAGE MODEL

903 In this paper, the term Large Vision-Language Models (LVLMs) refers to deep learning models de-
904 signed to process joint visual and textual data, built upon foundational LLMs. Specifically, LVLMs
905 integrate robust Large Language Models (LLMs) with pre-trained Vision encoders to enhance accu-
906 racy in understanding and generating language and vision-related content.

907 Typically, an LVLM is comprised of a vision encoder, a language encoder (i.e., an LLM), and a
908 cross-modal alignment network. The training process for LVLMs involves three primary stages.
909 Initially, the vision and language encoders undergo pre-training on extensive unimodal datasets,
910 focusing on image and text data separately. Subsequently, these encoders are aligned through pre-
911 training on image-text alignment, enabling the LLM to generate meaningful texts corresponding to
912 given images. Finally, the whole model undergoes further fine-tuning on image-text instructions,
913 enhancing its ability to provide satisfactory responses to natural language queries related to specific
914 images. Notably, during the second and third stages, selective fine-tuning of individual components
915 can be performed instead of conducting comprehensive parameter adjustments.

916 Once the visual encoder and the LLM are effectively aligned, the resulting LVLM exhibits superior
917 visual comprehension capabilities. It not only captures the visual semantics of objects within an
image but also delves into linguistic semantics by leveraging the parametric knowledge embedded

918 in the LLM, achieving enhanced performance across various vision language tasks, such as image
919 captioning.
920

921 E OBJECT HALLUCINATION AND CHAIR METRICS 922

923 **Object Hallucination:** In literature, the term "object hallucination" denotes a phenomenon
924 wherein a model generates descriptions or captions containing objects that are either inconsistent
925 with or entirely absent from the target image. Object hallucination can be understood and defined at
926 various semantic levels. At its simplest, it pertains to discrepancies at the object level, though more
927 nuanced interpretations may extend to the attributes or characteristics of objects. This study focuses
928 on object-level object hallucinations within model-generated captions, deferring finer-grained anal-
929 yses of object hallucinations—such as those related to quantity, attributes, and positions—to future
930 investigations.
931

932 **CHAIR:** The Caption Hallucination Assessment with Image Relevance (CHAIR) Rohrbach et al.
933 (2018a) stands as a widely recognized standard for gauging the occurrence of object hallucination in
934 image captioning tasks. This metric operates by scrutinizing the actual objects depicted in an image
935 and subsequently determining the percentage of referenced objects in the generated caption that do
936 not correspond to objects within the image itself. Two distinct variants of CHAIR are employed
937 to measure object hallucination: $CHAIR_s$, which evaluates object hallucination at the caption level,
938 and $CHAIR_o$, which assesses object hallucination at the object level. Mathematically, the metrics are
939 defined as follows:

$$940 \quad CHAIR_o = \frac{\# \{ \text{hallucinated objects} \}}{\# \{ \text{all objects in prediction} \}} \quad (12)$$

$$941 \quad CHAIR_s = \frac{\# \{ \text{captions with hallucinated objects} \}}{\# \{ \text{all captions} \}}. \quad (13)$$

942 F DESCRIPTION OF LVLM MODELS USED AS BASELINE 943

944 The evaluated LVLMs basically consist of three parts: a visual encoder, an alignment model, and a
945 large language model. All the above models have been tuned on collected visual instruction data
946

947 **mPlug-Owl** mPLUG-Owl Ye et al. (2023), is a novel training method that enhances LLMs with
948 multi-modal capabilities by integrating foundational LLM training, a visual knowledge module, and
949 a visual abstractor module. This approach supports various modalities and enhances both unimodal
950 and multimodal abilities through collaborative learning. mPLUG-Owl employs a two-stage training
951 process to align image and text data, leveraging LLM assistance while preserving and enhancing its
952 generative capacities. Initially, the visual knowledge and abstractor modules are trained using a fixed
953 LLM module to align image-text pairs. Subsequently, language-only and multi-modal supervised
954 datasets are utilized to fine-tune a Low-Rank Adaptation (LoRA) module on LLM and the abstractor
955 module while keeping the visual knowledge module frozen.
956

957 **LLaVA** uses a linear projector to map visual token as a soft-prompt into LLM input tokens. LLaVA
958 has a two-stage training, where the initial stage focuses on simple caption pretraining solely for the
959 linear projector, while the subsequent stage finetunes both the projector and LLM on instruction
960 data. Instruction data leverages language-only GPT-4 by inputting visual ground truth from COCO
961 dataset.
962

963 **InstructBLIP** adopts the BLIP-2 architecture, and is distinguished by its training of a Q-former,
964 which bridges the frozen vision encoder and LLM. InstructBLIP's instruction fine-tuning spans
965 across 26 distinct datasets.
966

967 G INSTRUCTION TEMPLATE FOR DETAILED IMAGE CAPTIONING IN COCO 968 DATASET 969

970 We use Instruction Templates to generate long, detailed captions. During training, the prompt is
971 randomly selected to query the LVLM. The Instruction Templates are at below:

- 972 • (Image)A detailed image caption:
- 973 • (Image)A detailed image description:
- 974 • (Image)Write a long description for the image.
- 975 • (Image)Describe the content of the image in detail.
- 976 • (Image)Can you explain clearly what you see in the image?
- 977 • (Image)Could you describe clearly what you perceive in the photo?
- 978 • (Image)Please provide a detailed depiction of the picture.
- 979 • (Image)Provide a detailed description of the given image.
- 980
- 981
- 982

983 H HARD PROMPT DESIGN

984 We have developed a set of "hard prompts" intended to be appended at the beginning of the input
 985 instruction, aiming to mitigate object hallucination in the model's generated captions. Each prompt
 986 is meticulously crafted to target specific sources of object hallucination, strategically guiding the
 987 model away from potential pitfalls. Below is the comprehensive list of prompts:

- 988 • Directly prohibit object hallucination : "Please don't hallucinate the objects in the image"
- 989 • Emphasize concrete details: "Provide captions based on specific, easily identifiable ele-
 990 ments in the image."
- 991 • Prioritize realism: "Generate captions that reflect plausible scenarios and avoid fantastical
 992 or improbable elements."
- 993 • Stick to visible entities: "Describe only what is clearly visible in the image and avoid
 994 making assumptions about hidden or obscured objects."
- 995 • Be conservative in interpretation: "Refrain from extrapolating beyond what is evident in
 996 the image; captions should stay closely tied to observable elements."
- 997 • Avoid creative interpretations: "Discourage the generation of captions that involve imagi-
 998 native or metaphorical representations of the scene."
- 999 • Limit descriptive scope: "Keep captions focused on the central objects or subjects in the
 1000 image, avoiding unnecessary details or peripheral elements."
- 1001 • Minimize speculative language: "Generate captions with certainty, avoiding speculative
 1002 language or uncertain descriptions of the depicted scene."
- 1003 • Resist contextual speculation: "Do not create captions that rely on external context or back-
 1004 ground information not present in the image."
- 1005 • Steer clear of abstract concepts: "Refrain from incorporating abstract or conceptual ideas
 1006 into the captions; stick to tangible, visible elements."
- 1007 • Encourage literal language: "Favor literal and straightforward language in captions, avoid-
 1008 ing figurative expressions or interpretations."
- 1009
- 1010
- 1011
- 1012

1013 I DETAILED ABOUT PROMPT TUNNING

1014 Image captioning with the Large Vision Language Model (LVLM) represents a crucial text genera-
 1015 tion task. Departing from the traditional classification approach, which assesses the probability of
 1016 an output class given input as $P(y|X, Z)$, where X comprises tokens representing the instruction,
 1017 y denotes a single class label, and Z contains tokens representing an image, we now adopt a condi-
 1018 tional generation perspective. In this paradigm, Y signifies a sequence of tokens that form a caption.
 1019 The captioning process by Large Vision Language Models is expressed as $P_\theta(Y|X, Z)$, where θ
 1020 represents the model's weights.

1021 Prompting involves augmenting the model's generation of Y by providing additional context for it
 1022 to rely on. This is achieved by prefixing a sequence of tokens, G , denoted as $\{g_1, g_2, \dots, g_k\}$, to
 1023 the input X , such that enabling the model to enhance the likelihood of generating the correct Y :
 1024 $P_\theta(Y|[G; X], Z)$. Throughout this process, the model parameters, θ , remain unchanged. Optimal
 1025

prompt selection can be achieved through manual exploration of prompt tokens, known as *Hard Prompting*, or by representing G with dedicated parameters, ϕ , which model the embeddings of these tokens. These parameters are then refined using gradient descent. This technique is termed *Soft Prompting*. Consequently, our updated conditional generation is expressed as $P_{\theta; \phi}(Y|[G; X], Z)$, and it can be trained by maximizing the reward through backpropagation, with gradient updates solely applied to ϕ .

The modeling of Soft Prompting is straightforward. When presented with a sequence of n tokens, $\{x_1, x_2, \dots, x_n\}$, the initial step undertaken by LVLM involves embedding these tokens to create a matrix $X_e \in \mathbb{R}^{n \times e}$, where e denotes the dimension of the embedding space. Our soft prompts are expressed as a parameter $G_e \in \mathbb{R}^{k \times e}$, with k being the length of the prompt. Subsequently, the prompt is concatenated to the embedded input, resulting in a unified matrix $[G_e; X_e] \in \mathbb{R}^{(k+n) \times e}$, which is then processed through the LVLM as per usual. During training, our models are designed to maximize the return of Y . However, it is noteworthy that only the prompt parameters G_e undergo updates, ensuring the model learns to effectively utilize the provided prompts while keeping other parameters fixed.

J DATASET DESCRIPTION

Visual Genome contains Visual Question Answering data in a multi-choice setting. It consists of 101,174 images from MSCOCO with 1.7 million QA pairs, 17 questions per image on average. Compared to the Visual Question Answering dataset, Visual Genome represents a more balanced distribution over 6 question types: What, Where, When, Who, Why and How. The Visual Genome dataset also presents 108K images with densely annotated objects, attributes and relationships.

The MS COCO (Microsoft Common Objects in Context) dataset is a large-scale object detection, segmentation, key-point detection, and captioning dataset. The dataset consists of 328K images.

Splits: The first version of MS COCO dataset was released in 2014. It contains 164K images split into training (83K), validation (41K) and test (41K) sets. In 2015 additional test set of 81K images was released, including all the previous test images and 40K new images.

Based on community feedback, in 2017 the training/validation split was changed from 83K/41K to 118K/5K. The new split uses the same images and annotations. The 2017 test set is a subset of 41K images of the 2015 test set. Additionally, the 2017 release contains a new unannotated dataset of 123K images.

The dataset has annotations for:

- object detection: bounding boxes and per-instance segmentation masks with 80 object categories.
- captioning: natural language descriptions of the images.
- keypoints detection: containing more than 200,000 images and 250,000 person instances labeled with keypoints (17 possible keypoints, such as left eye, nose, right hip, right ankle).
- stuff image segmentation: per-pixel segmentation masks with 91 stuff categories, such as grass, wall, sky.
- panoptic: full scene segmentation, with 80 thing categories (such as person, bicycle, elephant) and a subset of 91 stuff categories (grass, sky, road).
- dense pose: more than 39,000 images and 56,000 person instances labeled with DensePose annotations – each labeled person is annotated with an instance id and a mapping between image pixels that belong to that person body and a template 3D model. The annotations are publicly available only for training and validation images.

K QUANTITATIVE RESULT

We present quantitative results on the COCO dataset. The Base refers to the InstructBlip model, while Our model is the fine-tuned version of InstructBlip using our framework. Hallucinated objects are highlighted in bold and red for clarity.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133



Base: The image displays a collection of various vegetables and greens, all neatly arranged on top of an old-fashioned tablecloth. There are multiple rows with different types of produce: some have white ends while others possess more green coloration throughout the stalks or leaves. Various **knives** can be seen nearby in case they're needed for cutting up any larger portions during preparations later downstream from this setup stage involving setting out fresh veggies to drying/dehydrating them by placing those already cut into **glass jars** also known as "pickling"....

Our: A dining table displays a vibrant array of fresh vegetables, including both leafy greens and root veggies. Various colors like purple-tinged kale are present among the organic produce on display next to each other. In total, there are about 10 different kinds of plants arranged in an attractive composition...



Base: The image showcases a pair of zebras playfully fighting in the grass. They may be playing together rather than engaging in territorial battles commonly associated between wild animals such as **lions** vs **hyenas** which usually happen for food sources only.

Our: In this image, several zebras are interacting with each other in a lush green field. Some zebras are standing on top of each other, while others are running around the field. There are several trees visible in the background, adding to the natural setting.



Base: The scene depicts a busy urban street, with several cars parked near the houses on either side. One of these vehicles is situated in front and closer to you than others it appears dark brown or red colorized as "maroon". A group of **people** are congregated at one end; they appear more prominent towards left-center behind three **women** sitting together underneath an **umbrella** while two **men** walk past them holding different beverages nearby: **wine bottle** versus **coffee mug** respectively close by their hands.

Our: In this image, a red van is parked on the side of a road in front of a house. The van is situated near a stop sign, which is located close to the front of the house. There are several trees visible in the background, adding to the scenic setting.



Base: The image features a toilet bowl placed on the ground near an old **couch**. Next, there are several stacks of boxes and bags surrounding it. A **television** is situated nearby with its base visible towards left side next section by **stairs** leading downward likely used during restoring activities.

Our: In this image, there is a white toilet sitting on the ground in a junkyard. The toilet appears to be disassembled, with some of its parts scattered around the area.

L TRAINING SETTINGS FOR PROMPT FINE-TUNING

We present parameter settings in our framework during the training phase on the COCO dataset. Here we use 8 A6000 50G GPU for training. The information is shown in table below:

Parameters	
Max Soft Prompt Length n	20
Max Episode Length H	150
Mini-Batch Size	64
PPO Learning Rate	0.0002
Kullback–Leibler (KL) Initial Coefficient λ	0.1
KL Target Coefficient	0.05
Number of Training Epochs	50
Reward Weight for Meteor α	0.1
Reward Weight for BERTScore β	1
Vocabulary Size \mathcal{V}	32000

Table 2: Training parameters of our framework for the COCO dataset

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187



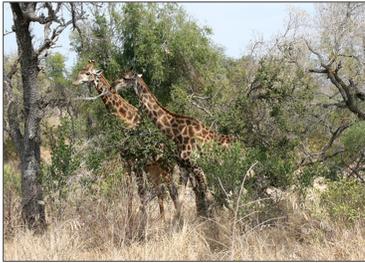
Base: The image depicts a small alleyway lined with buildings and balconies, featuring several **bikes** parked in the space between them. There are numerous **motorcycles** placed throughout this area. Additionally, there is another object that appears like it might be either one or two **cars** behind some nearby **trees** towards both sides near various parking spaces along city street setting.

Our: In this scenic view, a narrow cobblestone street winds its way through a charming old town. The street is lined with various buildings, some of which have multiple floors. There are several cars parked along the street, adding to the vibrant atmosphere.



Base: The image depicts a train moving along the railway tracks, with an open door on one of its cars. On either side there are multiple **smaller vehicles** parked nearby but do not appear involved into transportation as they resemble passenger automobiles without any visible loads travelers would board them at stops like passengers embark trains around station entrances while loading/unloading goods carried onto flatbed utility wagons (with wheels).

Our: At the train station, a large yellow train is parked on the tracks, waiting for its next departure. The train is decorated with an American flag, adding a patriotic touch to the scene.



Base: The image captures a pair of giraffes in their natural habitat, roaming across the forest. They are standing and walking among bushy trees covered with various leaves surround them like an umbrella canopy covering parts of this jungle landscape scene near where they live. There is at least three **people** present throughout the environment: one **person** sitting right next to another close by; several others scattered around both left-side areas.

Our: In this scenic image, two giraffes are standing in a lush, green forest. They appear to be socializing with each other while browsing through the trees. The giraffes are spread out across the forest, with one giraffe standing closer to the left side of the image, and the other further to the right.



Base: The image features a large pizza with various toppings in an open box. Two empty **bowls** can be seen near or underneath one side of another tabletop located offscreen towards left behind them. A **person** stands next to these tables having dinner together at some point during their meal preparation session where they enjoyed eating both savory pasta dishes as well as indulging themselves.

Our: A delicious mushroom pizza is sitting in a box, ready to be enjoyed. The pizza is sliced into multiple pieces, making it easy to share or enjoy on its own.



Base: The image depicts a group of young people gathered together in an old-fashioned setting. They are wearing various clothing and ties, with some sporting bowties or traditional dresses. A room is visible behind them as well featuring several **tables** set up inside it where guests may have been seated. Furthermore, there appears two **clock radios** placed towards different corners within that space giving evidence of activity happening here throughout time frames.

Our: In this vintage photograph, a large group of people are gathered together in a room. They are all dressed in suits and ties, creating a formal atmosphere. The group is engaged in socializing and interacting with each other.



Base: The image captures a beautiful day with trees in the background and **people** walking around. There is an area of dirt on one side near some water where three individuals can be seen engaging together or resting underneath two **umbrellas**. Several park **benches** are present along to each other as well towards both sides closer by the ground surface itself containing four seating places among them at least six **chairs** have been positioned further off but still easily visible through such foliage.

Our: In this scenic image, there is a red stop sign standing tall in the middle of a dirt road next to a fence. The sign is placed near a motorcycle, which can be seen in the background.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241



Base: The image showcases a shiny silver moped, parked neatly inside of an underground garage. There are two rear wheels on either side that make up most part of this compact vehicle's frame area near its tail section. A few **people** can be seen walking within various areas throughout the scene - specifically between right middle (one **person**), top centralized portion just beyond three riders sitting there beside another standing individual present alongside several **vehicles** also situated across four main locations.

Our: In this image, a sleek and modern motor scooter is parked in front of a brick wall. The motor scooter is silver in color and appears to be well-maintained. There are several motorcycles visible in the scene, creating a vibrant and lively atmosphere.