

# Mitigating Open-Vocabulary Caption Hallucinations

Anonymous ACL submission

## Abstract

While recent years have seen rapid progress in image-conditioned text generation, image captioning still suffers from the fundamental issue of hallucinations, namely, the generation of spurious details that cannot be inferred from the given image. Existing methods largely use closed-vocabulary object lists to mitigate or evaluate hallucinations in image captioning, ignoring most types of hallucinations that occur in practice. To this end, we propose a framework for addressing hallucinations in image captioning in the open-vocabulary setting, including quantifying their presence and optimizing to mitigate such hallucinations. Our *OpenCHAIR* benchmark leverages generative foundation models to evaluate open-vocabulary caption hallucinations, surpassing the popular CHAIR benchmark in both diversity and accuracy. To mitigate open-vocabulary hallucinations at the sequence level, we propose *MOCHa*, an approach harnessing advancements in reinforcement learning. Our multi-objective reward function explicitly targets the trade-off between fidelity and adequacy in generations without requiring any strong supervision. *MOCHa* improves a large variety of image captioning models, as captured by our *OpenCHAIR* benchmark and other existing metrics. We will release our code and models.

## 1 Introduction

Image captioning, the task of generating text that describes an image, is one of the most fundamental machine learning tasks combining vision and language. Unfortunately, *hallucinations* plague the current state-of-the-art (SOTA), making it less usable for practical tasks that require confidence in the factual correctness of generated captions. Consider, for instance, the images in Figure 1. SOTA image captioning models can generate text that is highly semantically related to its associated imagery, but also contains spurious details (“*dimly*

BLIP

+*MOCHa*



*Dimly shining* coffee. A glass mug of  
foam drink on top of a wooden table  
with a brown donut

Figure 1: Hallucinated details (shown as *highlighted text*) are prevalent in the outputs of modern image captioning models, such as the above generation sampled from BLIP (Li et al., 2022a). Considering hallucinations in the open-vocabulary setting, our RL-based *MOCHa* framework optimizes captioning models to output detailed, valid captions while avoiding such hallucinations, as illustrated in the right column (+*MOCHa*).

*shining*”, “*brown donut*”). Such hallucinated spurious details either damage user confidence or lead to uncritical acceptance of fallacious (and even potentially dangerous) generated content (Chong et al., 2022; McGowan et al., 2023; Chong et al., 2023).

Hallucinations may take a variety of forms in text, including complex multi-word phrases of various syntactic roles. However, prior work addressing hallucinations in image captioning has largely focused on detecting or mitigating hallucinations by using closed-vocabulary object lists. While this simplifies the problem under consideration, it fails to capture the diversity of hallucinations observed in modern image captioning models. Thus, we propose a framework for both quantifying and mitigating hallucinations in the open-vocabulary setting.

While established benchmarks and metrics for quantifying hallucinations in captioning models exist for closed-vocabulary object sets, they do not exist (to our knowledge) in an open-vocabulary setup. We introduce *OpenCHAIR*, a new benchmark for quantifying hallucinations in an open-vocabulary setting. By generating data and performing evaluation with text-to-image models and large language models (LLMs), we can capture and accurately quantify a wide variety of hallucination types with-

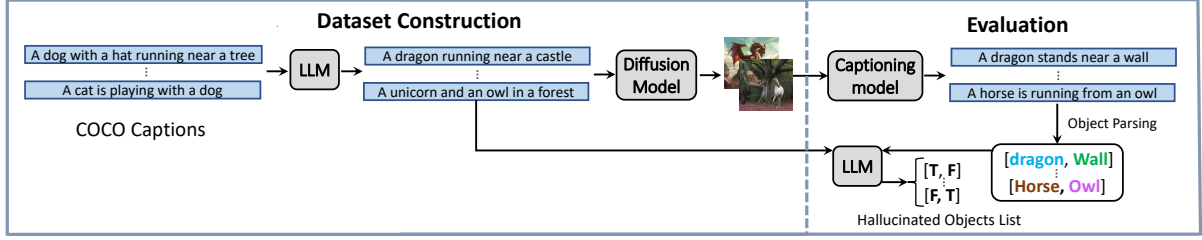


Figure 2: **The *OpenCHAIR* Benchmark.** We illustrate the construction of the *OpenCHAIR* benchmark via an LLM and text-to-image generation model, and its usage for evaluating image captioning models. We first use captions from MS-COCO as seeds to generate diverse synthetic captions. Using syntactic parsing and filtering heuristics, we select for captions containing various open-vocabulary objects. We then generate images corresponding to these captions, producing our benchmark of images linked with object annotations. To evaluate a captioning model, we run it on this benchmark and compare predicted and GT object categories.

out being limited to a fixed set of categories. Our evaluations show that this outperforms the CHAIR closed-vocabulary metric (Rohrbach et al., 2018).

Having the metric, we turn to hallucination mitigation. A major cause for hallucinations in image captioning (and text generation in general) models stem from deficiencies in the standard language modeling (LM) objective. The *token-level* likelihood maximization LM objective does not directly optimize the *sequence-level* quality of generated text, and *factual groundedness* is inherently a sequence-level property of text. Yet, many prior works that directly optimize hallucinations in image captioning avoid the global sequence-level nature of hallucination by limiting their scope to a fixed set of possible object tokens, e.g. objects in MS-COCO (Biten et al., 2021; Liu et al., 2022; Petryk et al., 2023), which is incompatible with an open-vocabulary setting.

To mitigate hallucinations in an open-vocabulary setting, we introduce *MOCHA*, a Multi-Objective reinforcement learning (RL) based approach for Mitigating Open-vocabulary Caption Hallucinations. We observe that RL applied to caption fidelity alone fails to preserve the semantic adequacy (i.e. descriptiveness) of output text, while optimizing for the latter does not enforce factually grounded text. Our key insight is that these two goals can be jointly optimized at the sequence-level by applying RL with a multi-objective reward function. Furthermore, we perform this optimization fully automatically by leveraging SOTA text-based learned metrics, without requiring direct supervision. By considering hallucinations in an open setting, we are able to improve performance across diverse hallucination types, as demonstrated by our *OpenCHAIR* benchmark as well as other metrics. Moreover, we

show that our approach can be flexibly applied to a variety of captioning architectures and sizes.

Explicitly stated, our key contributions are: (i) *OpenCHAIR*, a benchmark for open-vocabulary hallucinations in image captioning. (ii) *MOCHA*, a framework for optimizing a wide array of VLMs to produce high-quality factually-grounded output. (iii) Experiments showing the advantage of *OpenCHAIR* for measuring hallucinations in the open setting, and of *MOCHA* for reducing them.

## 2 The *OpenCHAIR* Benchmark

To measure object hallucination in the open-vocabulary settings, we propose the *OpenCHAIR* (OCH) benchmark. *OpenCHAIR* modifies the previous object hallucination metric CHAIR (Rohrbach et al., 2018), by relaxing its strong reliance on the object annotations in the MS-COCO dataset, which constitute only 80 object types. We provide an overview of *OpenCHAIR* below; further details on the construction and contents of the dataset, prompts used, and other implementation details are provided in the appendix.

In order to create a new benchmark that enables measuring the hallucination rate of arbitrary objects, while still maintaining high quality ground-truth captions, we use the pipeline illustrated in Figure 2. We first prompt the LLM Llama-2 (Touvron et al., 2023) with few-shot examples of image captions from MS-COCO, having it generate captions with a similar style but containing diverse details (and in particular, objects that are likely not contained in the closed set of MS-COCO object labels). We then parse these synthetic captions with a syntactic parsing model, identify nouns with high concreteness scores (Brysbaert et al., 2014) (as these generally represent concrete objects), and balance the generated captions among object types



“A green emerald is perched on a rock in a cave.” “A group of mushrooms in the forest.” “A dog dressed as a human with a wig and eyeglasses.”

Figure 3: **OpenCHAIR Examples.** We show examples of images from the *OpenCHAIR* benchmark along with their accompanying ground-truth captions, illustrating its diverse coverage of object types. Long captions are truncated due to space considerations.

to cover a wide array of objects. Subsequently, we utilize the text-to-image diffusion model Stable Diffusion XL (Podell et al., 2023) to generate an image from these newly formed caption. This process results in a dataset that consists of synthetic images with corresponding captions including diverse, open-vocabulary objects. Figure 3 shows examples of image-captions pairs from *OpenCHAIR*.

We evaluate captioning models on *OpenCHAIR* as follows: After generating captions for each image in the *OpenCHAIR* dataset, we parse them to identify objects as described above. For each extracted object  $o$ , we compare it to the ground-truth synthetic caption  $c$  by prompting an LLM, asking it whether an image with caption  $c$  contains the object  $o$  and using its answers to count hallucinations. Following CHAIR, we calculate the hallucination rate as  $n_h/n_{tot}$ , where  $n_h$  is the number of hallucinated objects (yes answers) and  $n_{tot}$  is the total number of objects considered. Figure 4 illustrates the difference between *OpenCHAIR* evaluation and the closed-vocabulary CHAIR metric.

### 3 The MOCHA Framework

To mitigate captioning hallucinations in the open-vocabulary setting, we propose *MOCHA*, an RL-based pipeline using SOTA methods for stable reinforcement along with a carefully designed reward function that jointly optimizes for caption fidelity and semantic adequacy. Figure 5 presents it. We turn to describe the learning procedure and objectives used in *MOCHA*. We start with preliminaries, then describe the reward function that *MOCHA* optimizes (Section 3.1), and finally present the RL algorithm used for optimization (Section 3.2).

**Preliminaries.** In general, RL views a model as an *agent* that interacts with the external *environment* and receives a *reward*, learning to optimize for this

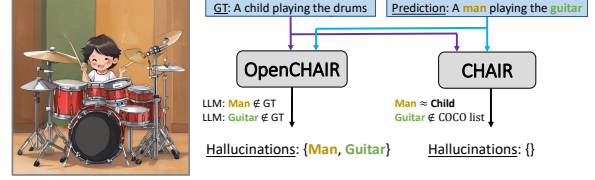


Figure 4: **OpenCHAIR vs. CHAIR.** In the above the predicted object *guitar* would not be counted by CHAIR since it is not in its fixed vocabulary, while *man* would not be classified as a hallucination since it is defined by CHAIR as a synonym of *child*. In contrast, *OpenCHAIR*’s LLM classifies both as hallucinations.

reward via exploring the environment (Sutton and Barto, 2018). In the case of image captioning, this model is a VLM operating in an environment of images and reference captions (Rennie et al., 2017). During training, the agent generates a caption by sampling from its own predicted distribution as shown in Figure 5 (left), receiving a reward based on an estimate of the caption quality. After collecting a full batch of rewards, a RL optimization step is applied as shown in Figure 5 (right), and this process repeats iteratively until convergence.

We use the following notation: Let  $T$  and  $I$  be the sets of possible texts and images, with joint distribution  $X$ . Given image  $i \in I$ , an image captioning model  $M$  with weights  $\theta$  induces a conditional probability distribution  $\pi_\theta(\cdot|\cdot)$  over generated captions  $\hat{c} \in T$  conditioned on images  $i \in I$ . In the RL context, we refer to  $\pi_\theta$  as the *policy*. A *reward function*  $r : T \times T \times I \rightarrow \mathbb{R}$  assigns *reward* (or score)  $r(\hat{c}; c, i)$  to generated caption  $\hat{c}$  relative to ground-truth caption  $c$  and image  $i$ .

#### 3.1 Reward Function

We wish to optimize for the competing objectives of output fidelity (low hallucination rate) and adequacy (including sufficient details to describe the input image), as optimizing for one of these alone causes the other to deteriorate (as shown in our ablations). We also wish to preserve other desired generation properties such as fluency and diversity. To achieve this, we design a reward function combining multiple objectives as follows:

**Fidelity Objective.** ( $r_f$ ). In order to measure output fidelity to the input image, we use the GT reference captions as a proxy for comparison, checking for logical consistency with generated caption via a pretrained Natural Language Inference (NLI) model. This outputs the probability  $\bar{p}(\hat{c}, c)$  that the generated text  $\hat{c}$  logically contradicts  $c$ , serving as a strong signal for fidelity, as details which con-



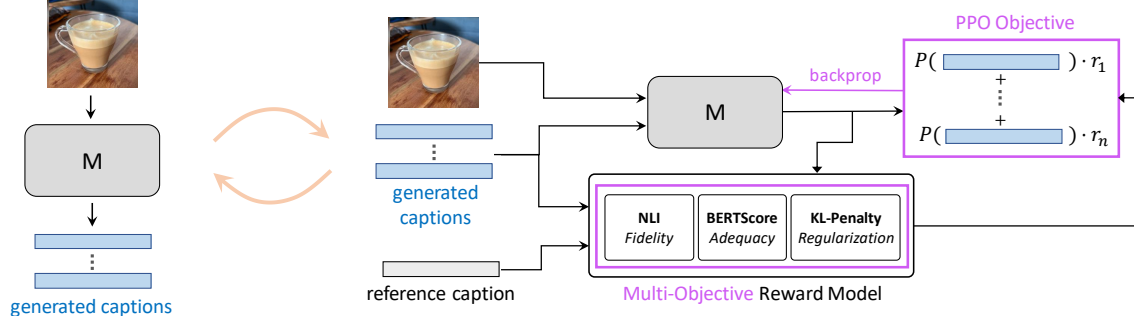


Figure 5: **MOCHA scheme**. The algorithm iteratively collects a minibatch of data from an image captioning model  $M$  (left side) and then applies an optimization step to the captioning model (right side). The multi-objective reward reinforces  $M$  to produce captions closer to the high-scoring captions and further from the low-scoring captions.

tradict ground-truth information about the image are guaranteed to be hallucinations. We scale to the range  $[-1, 1]$  by using  $r_f(\hat{c}; c) := 1 - 2\bar{p}(\hat{c}, c)$  as the fidelity reward. We implement this with BART (Lewis et al., 2019) fine-tuned on the MNLI dataset (Williams et al., 2018). We average values over all reference captions.

**Adequacy Objective.** ( $r_a$ ). To measure adequacy (whether the output caption contains sufficient detail), we use BERTScore (Zhang et al., 2019), a pretrained model measuring text quality relative to ground-truth references. We calculate its F1 value, scaled scale to be approximately in the range  $[-1, 1]$  as described in the appendix.

**KL Regularization.** Following prior work (Jaques et al., 2017, 2019; Ziegler et al., 2020; Stiennon et al., 2020; Ouyang et al., 2022), we add a Kullback–Leibler (KL) divergence penalty to the reward model which constrains the agent to stay close to its initial policy  $\pi_0$ . This serves to prevent mode collapse (i.e. preserving diversity of outputs) and adversarial policies which over-optimize the reward function. The KL penalty adds a term proportional to  $K(\hat{c}; i) := -\log(\pi_\theta(\hat{c}|i)/\pi_0(\hat{c}|i))$  to the reward, which limits the agent from excessively distancing itself from the initial policy.

**Combined Objective.** Our total reward function takes the form  $r(\hat{c}; c, i) := \alpha \cdot r_f(\hat{c}; c) + (1 - \alpha) \cdot r_a(\hat{c}; c) + \beta K(\hat{c}; i)$ , where  $\alpha$  and  $\beta$  are positive scalars controlling the trade-off between objectives.

### 3.2 Learning Procedure

To optimize for caption generations that satisfy the desired properties (described above in Section 3.1), we adopt the Proximal Policy Optimization (PPO) RL algorithm (Schulman et al., 2017), which has been used by recent works on text generation as discussed in Section 5. This is a *policy gradient* algorithm, meaning that it optimizes the parameters  $\theta$

in order to (approximately) maximize the expected reward  $L(\theta) = E_{i, c \sim X, \hat{c} \sim \pi_\theta(\hat{c}|i)} [r(\hat{c}; c, i)]$ . PPO extends the REINFORCE algorithm (Sutton and Barto, 2018), also known as SCST in the context of image captioning (Rennie et al., 2017), by using a clipped surrogate objective to avoid instabilities.

## 4 Experiments and Results

### 4.1 OpenCHAIR Analysis and Comparison to CHAIR

We analyze the utility of *OpenCHAIR* by comparing its distribution of objects to the existing closed-vocabulary CHAIR metric, as well as by performing a human evaluation to compare their correlations to human judgements of hallucinations.

In the first column of Table 1 and in Figure 14 (appendix), we show the difference in the number of unique object types found in CHAIR and *OpenCHAIR*. The open-vocabulary design of *OpenCHAIR* enables a significantly larger coverage of object types, more than ten times as many as used to evaluate the CHAIR metric. This is also reflected qualitatively, as the closed-vocabulary benchmark is missing many common object types, including daily objects like *shoe* and *guitar* (see the left image in Figure 6 for a visual example). In contrast, our benchmark includes diverse object types, such as: *pearl*, *tiger*, *sand*, *tricycle*, *corkscrew*, *toy*, *charcoal*, *text*, *pine-cone*, *grandfather*, *chocolate*, *wheelchair*, *wand*, etc. A list of all additional objects (those not included in CHAIR) can be found in the accompanying file `openchair_objects.txt`. Another source of confusion in CHAIR is its synonym list. See Figure 4 and the discussion below.

We show that *OpenCHAIR* evaluations are grounded in human intuitions via a manual evaluation, comparing its performance to that of CHAIR. For each benchmark (*OpenCHAIR* and CHAIR), we generate captions for a random subset of its



	# Obj Types	Balanced Accuracy			
		BLIP2	BLIP-L	GIT-B	OFA-L
CH	80*	0.844	0.774	0.899	0.810
OCH	<b>1360</b>	<b>0.945</b>	<b>0.944</b>	<b>0.943</b>	<b>0.930</b>

Table 1: **Human Evaluation of *OpenCHAIR* and CHAIR.** We perform a manual evaluation of *OpenCHAIR* and CHAIR object-level predictions, as described in Section 4.1. As seen above, *OpenCHAIR* covers a much larger variety of unique object types while also outperforming CHAIR in per-object predictive accuracy (of whether the given object is present or hallucinated). \*CHAIR includes also a synonym list.

dataset and manually check object-level decisions (predicted as existing or hallucinated) for over 400 random objects. Results using various captioning models are found in Table 1. As the presence of hallucinations is highly imbalanced (the large majority of predicted objects are not hallucinated), we report balanced accuracy. We provide further details in the appendix, including full confusion matrices.

Surprisingly, although operating over a much more diverse scope, *OpenCHAIR* achieves higher accuracy than CHAIR. We identify that this issue stems from CHAIR’s heavy reliance on coarse synonym lists, as seen in Figure 6 (right). By assessing whether pairs of object names match using a knowledgeable LLM, *OpenCHAIR* performs finer-grained hallucination measurements and achieves superior accuracy even in the more general open-vocabulary setting. We note that this reflects a trade-off between true and false positives, as predicted objects may not be found in *OpenCHAIR* ground-truth lists despite being present in the accompanying images, due to the limited descriptive capacity of text used to generate images. See more details in the Appendix (Tables 3 and 4).

## 4.2 *MOCHA* Implementation Details

We test image captioning with *MOCHA* on various SOTA image captioning models of varying architectures and across various sizes. In particular, we test BLIP (Li et al., 2022a), BLIP-2 (Li et al., 2023a) and GIT (Wang et al., 2022). Following standard practice in RL-based image captioning, we use models that have first been fine-tuned on with a standard language modeling loss on the captioning dataset, and then applying PPO reinforcement with our reward function ( $\alpha = 0.5$ ). See the appendix for model checkpoints, parameter counts, and further training settings and hyperparameters.

We test our method on the MS-COCO (Lin et al.,

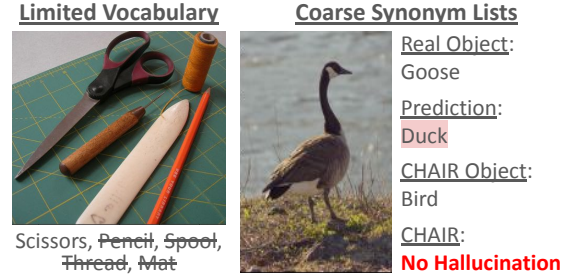


Figure 6: **CHAIR Limitations.** The left image exhibits CHAIR’s limited vocabulary. Out of all objects predicted by BLIP2, *Scissors* is the only object CHAIR considers during the evaluation. The right image illustrates a limitation stemming from CHAIR’s use of a fixed list of synonyms to coarsely aggregate different, semantically similar objects. Hallucinations that occur within the same synonym group are considered as a correct detection; in this example both *Goose* and *Duck* are defined as synonyms of *Bird* even though the image does not display a duck (but rather a goose).

Model	Quality		Hallucination			
	B@4↑	C↑	CH <sub>i</sub> ↓	CH <sub>s</sub> ↓	OCH ↓	$\bar{p}$ ↓
BLIP	41.5	138.4	2.3	3.5	16.4	0.244
BLIP+L	5.5	0.0	12.1	35.4	28.63	0.321
BLIP+T	41.3	137.4	<b>1.9</b>	<b>2.8</b>	16.4	0.241
BLIP+M	<b>41.9</b>	<b>139.6</b>	2.1	3.1	<b>15.4</b>	<b>0.206</b>

Table 2: **Comparison To Prior Works.** Measured over MS-COCO for BLIP-Large. +L/T/M refer to LURE, TLC-A, and *MOCHA* respectively. BSc and  $\bar{p}$  denote BERTScore and NLI contradiction probability rewards. B@4, C, CH, OCH, and  $\bar{p}$  denote BLEU-4, CIDEr, CHAIR, OpenCHAIR, and NLI  $p(\text{contr.})$  metrics respectively. Best results are shown in **bold**.

2015) captioning benchmark, using the data split of Karpathy and Fei-Fei (Karpathy and Fei-Fei, 2015) (113K items for training, 5K for evaluation). We report standard captioning metrics along with CHAIR (Rohrbach et al., 2018) and *OpenCHAIR* over generated captions (beam search decoding with 5 beams). We also provide NLI ( $\bar{p}$ ) and BERTScore values, directly optimized by *MOCHA*, as described in Section 3.1. In the appendix, we provide results on additional captioning datasets and metrics to further demonstrate generalization.

## 4.3 *MOCHA* Results

See Figure 7 for quantitative results of image captioning models on MS-COCO, where we show the relative improvement of optimizing the baseline SOTA captioning models with *MOCHA*. As is seen there, *MOCHA* improves measures of hallucina-

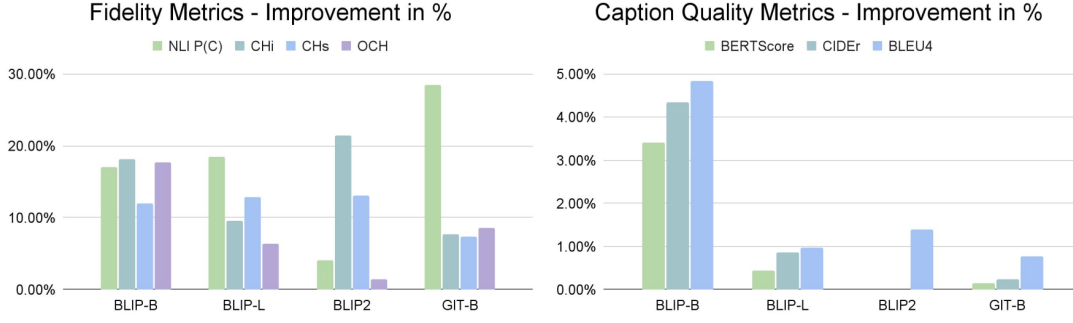


Figure 7: **Reducing Hallucinations While Maintaining Caption Quality.** We show the relative improvement of state-of-the-art VLM models when optimized using *MOCHA* optimization on the COCO Caption Karpathy test set. CH and OCH refer to Chair and *OpenCHAIR* respectively. All results are generated by using their officially provided checkpoints and hyperparameters. Full numeric results are provided in the appendix.

<b>B</b>	A man in a suit and tie standing by another man in a suit and tie	A person taking a tray of apples out of an oven	A man sitting on a couch talking on a cell phone
<b>B+M</b>	A man in a military uniform talking to a man in a suit and tie	A person taking a pan of food out of an oven	A man sitting on a couch using a laptop computer

Figure 8: **Qualitative results** of *MOCHA* applied to an image captioning model (BLIP-Large), along with baseline results without optimization (noted as B+M, B, respectively). We show captions (over COCO) produced from each model using beam search decoding with five beams. Hallucinated details are highlighted. The results illustrate that *MOCHA* encourages captions with high fidelity to the input image (avoiding hallucinations), while preserving a satisfying level of detail.

tions in image captioning while preserving or even enhancing standard measures of caption quality. We note that this is despite the fact that the trade-off between these qualities may degrade one or the other when using a sub-optimal reward weighting (as shown in Section 4.5 and observed in numerous prior works). We also provide qualitative examples in Figure 8, illustrating that the *MOCHA*-optimized model generates captions consistent with the input images while preserving a satisfying level of detail, consistent with our numeric results.

Our quantitative results show that *MOCHA* improves performance over base captioning models by most measures, across model architectures and sizes. This effect is seen not only among metrics that we directly optimize (NLI, BERTScore) but also among non-optimized metrics, measuring general caption quality (e.g. CIDEr), closed-

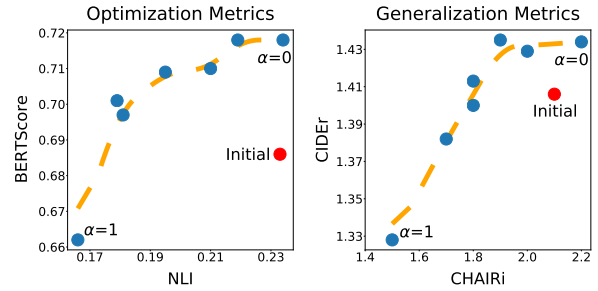


Figure 9: **Fidelity-Adequacy graphs** for pretrained (“initial”) and *MOCHA*-optimized BLIP models. As seen above, varying the reward weighting  $\alpha$  adjusts the trade-off between caption fidelity (x-axis) and adequacy (y-axis), with intermediate values outperforming the initial model (“Initial”). This holds both for metrics we directly optimize (left) and additional metrics (right), illustrating the generalization ability of our approach.

vocabulary hallucinations (CHAIR) and open-vocabulary hallucinations (*OpenCHAIR*). Along with our qualitative observations, this justifies our holistic approach to reducing hallucinations without restriction to a closed objects list.

#### 4.4 *MOCHA* Comparisons

In Table 2 we compare *MOCHA* to LURE (Zhou et al., 2024) and TLC-A (Petryk et al., 2023), current SOTA methods addressing VLM hallucinations, applied to the same pretrained BLIP model. LURE fails in the pure image captioning setting as its training procedure encourages long-form, highly detailed outputs. While these are in-distribution for instruction-tuned VLMs, they represent an increase in hallucinations relative to concise captions, as well as an extreme deviation from the reference texts; thus it degrades performance across metrics when applied to a captioning model such as BLIP. Regarding TLC-A, as it targets the objects in the closed-vocabulary object list of CHAIR, it shows an expected advantage in this metric, but degrades *OpenCHAIR* and other measures of caption quality,

contrasting with the overall improvement shown by our method. More details and results are provided in Appendix B.3, B.4 and C.4.

A number of prior works have proposed dedicated methods for reduced-hallucination image captioning, often using data modification or building multi-component pipelines applied to older vision-language backbones. In Table 8 (appendix), we provide a comparison between these methods and SOTA foundation VLMs applied as-is, reproducing results for the dedicated methods UD-L (Biten et al., 2021), CIIC (Liu et al., 2022), and COSNET (Li et al., 2022b). We find SOTA VLMs outperform these methods across all metrics, motivating our focus on optimization applied on top of modern foundation models.

#### 4.5 Ablations

We ablate the components of our reward function, finding that optimizing for fidelity alone degrades general caption quality, while optimizing for adequacy alone fails to improve hallucinations. This is seen in Figure 9 where extreme values of  $\alpha$  (0 or 1) correspond to the edges of the curves. Adjusting the parameter  $\alpha$  controlling the trade-off between objectives traces a *Pareto frontier* which outperforms the base model, showing that joint optimization of these objectives has a synergistic effect. The effects of each reward function component are also illustrated qualitatively in Figure 15 (appendix); removing  $r_f$  from the reward function leads to increased hallucinations, and removing  $r_a$  leads to captions that do not contain sufficient details. We provide full numeric results in the appendix, as well as ablating the effect of our chosen RL algorithm and of the KL-Penalty in our reward.

### 5 Related Work

**Measuring VLM Hallucinations.** A number of methods for measuring hallucinations in generated text have been proposed (Ji et al., 2023). In particular, various methods quantify hallucinations in the context of image-conditioned text generation, as summarized in Figure 10 (left). Among metrics that quantify hallucinations in predicted captions (“Prediction Assessing” in the figure), the existing CHAIR metric Rohrbach et al. (2018) explicitly quantifies object hallucinations by comparing tokens occurring in predicted captions to ground-truth object annotations. This requires a dataset

such as COCO that contains object annotations along with images, and assumes a fixed vocabulary of object identities. In our work, we demonstrate that this approach can be modified by leveraging advancements in LLMs and text-to-image generation models; our *OpenCHAIR* benchmark thus provides an explicit measure of open-vocabulary hallucinations in predicted captions via diverse ground-truth object annotations paired with generated images.

A handful of works have proposed more holistic measures of the fidelity of generated text with respect to an input image (the “Similarity Based” metrics of Figure 10) using embedding similarities or learned scores. CLIPScore (Hessel et al., 2022) propose CLIP cross-modal similarity for detecting mismatches between text and images, including hallucinations, and Shi et al. (2022) propose a similar embedding-based metric for video captioning. However, Xu et al. (2023) find that CLIP tends to assign high similarity to texts with minor modifications (“hard negatives”) that contradict the corresponding image. The Egoshots Semantic Fidelity metric (Agarwal et al., 2020) and VIFIDEL (Madhyastha et al., 2019) use embedding similarity between object annotations or detections in images and items in predicted captions. FAIEr (Wang et al., 2021) proposes a learned fidelity metric, which must be trained on automatically-generated scene graphs. While these metrics correlate with the presence of hallucinations, they are less interpretable as they do not provide a discrete count of hallucinations in a predicted caption.

Li et al. (2023b) propose POPE, which compares a list of ground-truth objects to the model’s answers when asked if each object is present. While this is open-vocabulary, it differs from our setting as it does not score predicted captions but rather assesses a VQA model’s general knowledge (indicated as “Model Assessing” in Figure 10).

**Reducing VLM Hallucinations.** Various methods for mitigating hallucinations in image captioning have been proposed, as illustrated in Figure 10 (right). Until recently, research on mitigating hallucinations in captions has largely considered object (noun) hallucinations, typically confined to a closed vocabulary (e.g. objects defined in MSCOCO). UD-L (Biten et al., 2021) identifies object hallucinations with bias towards the prior distribution of objects in context found in the training data, and proposes the use of synthetically debiased captions. CIIC (Liu et al., 2022) focuses on captioning models with a closed-vocabulary object



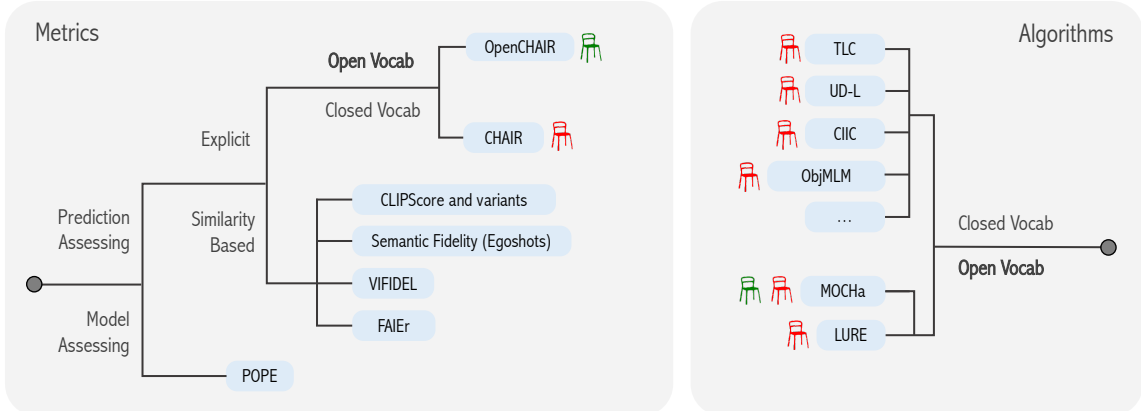


Figure 10: **VLM Caption Hallucination Taxonomy.** We illustrate metrics (left) and algorithms (right) for quantifying and mitigating hallucinations in image-conditioned text generation. We propose an explicit metric for measuring open-vocabulary hallucinations (*OpenCHAIR*) and an open-vocabulary hallucination mitigation algorithm (*MOCHa*). We mark each algorithm with the automatic hallucination rate metric with which it is evaluated (Green – *OpenCHAIR*, Red – *CHAIR*). Further details are provided in Section 5.

detection backbone, inserting components into the object detector and text decoder to reduce spurious correlations. TLC (Petryk et al., 2023) proposes a text decoding method applied to existing captioning models, to avoid generating COCO object tokens if they have insufficient confidence. Yin et al. (2023) combine closed-vocabulary object detection with LLM-guided decoding to avoid hallucinations in generated text. The more recent work ObjMLM (Dai et al., 2023) proposes masking objects from closed vocabulary lists as a training objective. Unlike these works, we mitigate hallucinations in the more challenging open-vocabulary setting. The contemporary work LURE (Zhou et al., 2024) proposes a method for the open setting, but their proposed approach (complementary to ours) was not evaluated in an open vocabulary setting due to the lack of an existing benchmark. Figure 10 illustrates which explicit hallucination metric was used to evaluate each algorithm.

As instruction-following VLMs rapidly develop, multiple concurrent works have considered hallucinations in related tasks such as visual question-answering (VQA), applying RL-based methods adopted from research on LLMs. These approaches train a reward model using a manually labelled dataset of hallucinations, then use this model for RL fine-tuning to reduce hallucinations in large VLMs (Gunjal et al., 2023; Sun et al., 2023a,b). These methods, which do not directly target our task, also require laborious human annotation to train a supervised reward model (while our approach does not require any explicit supervision). **Deep RL for VLM Text Generation.** Deep RL has been widely applied to text generation tasks.

One successful line of work optimizes such metrics for image captioning using an approach called Self-Critical Sequence Training (SCST) (Rennie et al., 2017; Stefanini et al., 2022). Another more recent development is the rise of deep RL for LLM alignment to user preferences. This commonly uses the Reinforcement Learning from Human Feedback (RLHF) framework, involving manual preference annotation followed by reinforcement-based optimization using a model to predict human preferences as a reward signal (Ziegler et al., 2020; Stiennon et al., 2020; Ouyang et al., 2022). Beyond LLMs, RLHF has been recently applied to aligning multimodal models with human preferences (Abramson et al., 2022). While such methods succeed in optimizing sequence-level properties, they often suffer from increased hallucinations as a side-effect of optimizing for human preferences or standard NLG sequence-level metrics (as illustrated in Appendix C.4).

## 6 Conclusion

We have shown the significance of operating in an open-vocabulary setting to effectively quantify and mitigate caption hallucinations. These are explicitly measured by our *OpenCHAIR* benchmark, and our *MOCHa* framework allows for optimizing captioning models to reduce such hallucinations while preserving caption quality. This reduction is demonstrated on our benchmark and other existing metrics. Our method and benchmark may be applied flexibly to a variety of model sizes and architectures, which we foresee providing a framework for future work on hallucination-aware image captioning.

## 7 Limitations

While the use of generative foundation models provides flexibility in evaluating open-vocabulary hallucination, it may inherit the limitations of these models including the lack of interpretability of their predictions. In addition, a potential limitation of our optimization approach is that it relies only on text despite the fact that it addresses the problem of image captioning that is fundamentally grounded in visual data. While our strategy achieves a consistent improvement across different models, the fact that it does not directly consider the image information may limit its performance.

We emphasize that our work does not solve the hallucination problem completely, although it presents a significant step towards this goal. Note also that we have focused in this work on the image captioning domain, while modern VLMs are often applied to diverse tasks such as VQA and visual instruction-following for which hallucinations also pose a significant challenge. We hope that our proposed strategy will pave the way for future research on hallucination reduction in all of these domains, in which open-vocabulary approaches also present significant promise.

## 8 Ethics Statement

This work focuses on measuring and mitigating hallucinations in visual-language models (VLMs). As such it is expected to increase the reliability of VLMs and the ability to measure their performance, which is important when using them in real world systems. This is expected to have a positive impact on the use of VLMs in the society. However, we do recognize that the foundation models used in the *OpenCHAIR* construction and evaluation pipeline and those used to calculate the *MOCHA* reward function could propagate biases. We anticipate further research into such biases before relying on our work beyond the research environment.

## References

- Josh Abramson, Arun Ahuja, Federico Carnevale, Petko Georgiev, Alex Goldin, Alden Hung, Jessica Landon, Jirka Lhotka, Timothy Lillicrap, Alistair Muldal, et al. 2022. Improving multimodal interactive agents with reinforcement learning from human feedback. *arXiv preprint arXiv:2211.11602*.
- Pranav Agarwal, Alejandro Betancourt, Vana Panagiotou, and Natalia Díaz-Rodríguez. 2020. Egoshots,

an ego-vision life-logging dataset and semantic fidelity metric to evaluate diversity in image captioning models. *arXiv preprint arXiv:2003.11743*.

Ali Furkan Biten, Lluís Gomez, and Dimosthenis Karatzas. 2021. [Let there be a clock on the beach: Reducing object hallucination in image captioning](#).

Marc Brysbaert, Amy Beth Warriner, and Victor Kuperman. 2014. Concreteness ratings for 40 thousand generally known english word lemmas. *Behavior research methods*, 46:904–911.

Leah Chong, Ayush Raina, Kosa Goucher-Lambert, Kenneth Kotovsky, and Jonathan Cagan. 2023. The evolution and impact of human confidence in artificial intelligence and in themselves on ai-assisted decision-making in design. *Journal of Mechanical Design*, 145(3):031401.

Leah Chong, Guanglu Zhang, Kosa Goucher-Lambert, Kenneth Kotovsky, and Jonathan Cagan. 2022. Human confidence in artificial intelligence and in themselves: The evolution and impact of confidence on adoption of ai advice. *Computers in Human Behavior*, 127:107018.

Wenliang Dai, Zihan Liu, Ziwei Ji, Dan Su, and Pascale Fung. 2023. Plausible may not be faithful: Probing object hallucination in vision-language pre-training. In *European Chapter of the Association for Computational Linguistics*, pages 2136–2148.

Anisha Gunjal, Jihan Yin, and Erhan Bas. 2023. Detecting and preventing hallucinations in large vision language models. *arXiv preprint arXiv:2308.06394*.

Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2022. [Clipscore: A reference-free evaluation metric for image captioning](#).

Jack Hessel, David Mimno, and Lillian Lee. 2018. Quantifying the visual concreteness of words and topics in multimodal datasets. *arXiv preprint arXiv:1804.06786*.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*.

Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#).

Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza, Noah Jones, Shixiang Gu, and Rosalind Picard. 2019. Way off-policy batch deep reinforcement learning of implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456*.

662	Natasha Jaques, Shixiang Gu, Dzmitry Bahdanau,	Alessia McGowan, Yunlai Gui, Matthew Dobbs, Sophia	717
663	José Miguel Hernández-Lobato, Richard E Turner,	Shuster, Matthew Cotter, Alexandria Selloni, Mar-	718
664	and Douglas Eck. 2017. Sequence tutor: Conserva-	ianne Goodman, Agrima Srivastava, Guillermo A	719
665	tive fine-tuning of sequence generation models with	Cecchi, and Cheryl M Corcoran. 2023. Chatgpt and	720
666	kl-control. In <i>International Conference on Machine</i>	bard exhibit spontaneous citation fabrication during	721
667	<i>Learning</i> , pages 1645–1654. PMLR.	psychiatry literature search. <i>Psychiatry Research</i> ,	722
		326:115334.	723
668	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	724
669	Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	725
670	Madotto, and Pascale Fung. 2023. Survey of halluci-	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	726
671	nation in natural language generation. <i>ACM Comput-</i>	2022. Training language models to follow instruc-	727
672	<i>ing Surveys</i> , 55(12):1–38.	tions with human feedback. <i>Advances in Neural</i>	728
673	Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-	<i>Information Processing Systems</i> , 35:27730–27744.	729
674	semantic alignments for generating image descrip-		
675	tions. In <i>Proceedings of the IEEE conference on</i>	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	730
676	<i>computer vision and pattern recognition</i> , pages 3128–	Jing Zhu. 2002. Bleu: a method for automatic evalu-	731
677	3137.	ation of machine translation. In <i>Proceedings of the</i>	732
		<i>40th annual meeting of the Association for Computa-</i>	733
678	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan	<i>tional Linguistics</i> , pages 311–318.	734
679	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,		
680	Veselin Stoyanov, and Luke Zettlemoyer. 2019.	Suzanne Petryk, Spencer Whitehead, Joseph E. Gon-	735
681	BART: denoising sequence-to-sequence pre-training	zalez, Trevor Darrell, Anna Rohrbach, and Marcus	736
682	for natural language generation, translation, and com-	Rohrbach. 2023. Simple token-level confidence im-	737
683	prehension. <i>CoRR</i> , abs/1910.13461.	proves caption correctness.	738
684	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi.	Dustin Podell, Zion English, Kyle Lacey, Andreas	739
685	2023a. Blip-2: Bootstrapping language-image pre-	Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna,	740
686	training with frozen image encoders and large lan-	and Robin Rombach. 2023. Sdxl: Improving latent	741
687	guage models. <i>arXiv preprint arXiv:2301.12597</i> .	diffusion models for high-resolution image synthesis.	742
		<i>arXiv preprint arXiv:2307.01952</i> .	743
688	Junnan Li, Dongxu Li, Caiming Xiong, and Steven	Steven J Rennie, Etienne Marcheret, Youssef Mroueh,	744
689	Hoi. 2022a. Blip: Bootstrapping language-image	Jerret Ross, and Vaibhava Goel. 2017. Self-critical	745
690	pre-training for unified vision-language understand-	sequence training for image captioning. In <i>Proceed-</i>	746
691	ing and generation. In <i>International Conference on</i>	<i>ings of the IEEE conference on computer vision and</i>	747
692	<i>Machine Learning</i> , pages 12888–12900. PMLR.	<i>pattern recognition</i> , pages 7008–7024.	748
693	Yehao Li, Yingwei Pan, Ting Yao, and Tao Mei. 2022b.	Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns,	749
694	Comprehending and ordering semantics for image	Trevor Darrell, and Kate Saenko. 2018. Ob-	750
695	captioning.	ject hallucination in image captioning. <i>CoRR</i> ,	751
		abs/1809.02156.	752
696	Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang,	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec	753
697	Wayne Xin Zhao, and Ji-Rong Wen. 2023b. Eval-	Radford, and Oleg Klimov. 2017. Proximal policy	754
698	uating object hallucination in large vision-language	optimization algorithms.	755
699	models.		
700	Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir	Yaya Shi, Xu Yang, Haiyang Xu, Chunfeng Yuan, Bing	756
701	Bourdev, Ross Girshick, James Hays, Pietro Perona,	Li, Weiming Hu, and Zheng-Jun Zha. 2022. Emscore:	757
702	Deva Ramanan, C. Lawrence Zitnick, and Piotr Dol-	Evaluating video captioning via coarse-grained and	758
703	lár. 2015. Microsoft coco: Common objects in con-	fine-grained embedding matching.	759
704	text.		
705	Bing Liu, Dong Wang, Xu Yang, Yong Zhou, Rui Yao,	Matteo Stefanini, Marcella Cornia, Lorenzo Baraldi,	760
706	Zhiwen Shao, and Jiaqi Zhao. 2022. Show, decon-	Silvia Cascianelli, Giuseppe Fiameni, and Rita Cuc-	761
707	found and tell: Image captioning with causal infer-	chiara. 2022. From show to tell: A survey on deep	762
708	ence. In <i>2022 IEEE/CVF Conference on Computer</i>	learning-based image captioning. <i>IEEE transac-</i>	763
709	<i>vision and Pattern Recognition (CVPR)</i> , pages 18020–	<i>tions on pattern analysis and machine intelligence</i> ,	764
710	18029.	45(1):539–559.	765
711	Pranava Madhyastha, Josiah Wang, and Lucia Specia.	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel	766
712	2019. VIFIDEL: Evaluating the visual fidelity of	Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,	767
713	image descriptions. In <i>Proceedings of the 57th An-</i>	Dario Amodei, and Paul F Christiano. 2020. Learn-	768
714	<i>annual Meeting of the Association for Computational</i>	ing to summarize with human feedback. <i>Advances</i>	769
715	<i>Linguistics</i> , pages 6539–6550, Florence, Italy. Asso-	<i>in Neural Information Processing Systems</i> , 33:3008–	770
716	ciation for Computational Linguistics.	3021.	771



772	Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu,	Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun	827
773	Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan	Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and	828
774	Gui, Yu-Xiong Wang, Yiming Yang, Kurt Keutzer,	Huaxiu Yao. 2024. Analyzing and mitigating object	829
775	and Trevor Darrell. 2023a. <a href="#">Aligning large multi-</a>	hallucination in large vision-language models. In	830
776	<a href="#">modal models with factually augmented rlhf.</a>	<i>ICLR.</i>	831
777	Zhiqing Sun, Yikang Shen, Hongxin Zhang, Qinhong	Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B.	832
778	Zhou, Zhenfang Chen, David Cox, Yiming Yang, and	Brown, Alec Radford, Dario Amodei, Paul Chris-	833
779	Chuang Gan. 2023b. <a href="#">Salmon: Self-alignment with</a>	tiano, and Geoffrey Irving. 2020. <a href="#">Fine-tuning lan-</a>	834
780	<a href="#">principle-following reward models.</a>	<a href="#">guage models from human preferences.</a>	835
781	Richard S. Sutton and Andrew G. Barto. 2018. <a href="#">Rein-</a>	<b>A Interactive Visualization</b>	836
782	<a href="#">forcement Learning: An Introduction</a> , second edition.		
783	The MIT Press.	For additional qualitative results, we refer the	837
784	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	reader to the interactive visualization tool provided	838
785	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	at the attached file <a href="#">index.html</a> . We provide im-	839
786	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	age captioning results using BLIP-Large with and	840
787	Bhosale, et al. 2023. Llama 2: Open founda-	without <i>MOCHA</i> for 350 randomly selected test	841
788	tion and fine-tuned chat models. <i>arXiv preprint</i>	images from MS-COCO (Lin et al., 2015) and	842
789	<i>arXiv:2307.09288.</i>	Flickr30K (Young et al., 2014).	843
790	Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie	To visually emphasize the hallucination rate in	844
791	Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and	the predictions, for each model we calculate the	845
792	Lijuan Wang. 2022. Git: A generative image-to-text	NLI contradiction probability <sup>1</sup> between the top	846
793	transformer for vision and language. <i>arXiv preprint</i>	beam and a ground-truth caption (which is depicted	847
794	<i>arXiv:2205.14100.</i>	below the image), and report the difference in the	848
795	Sijin Wang, Ziwei Yao, Ruiping Wang, Zhongqin Wu,	contradiction probability between the two models.	849
796	and Xilin Chen. 2021. Faier: Fidelity and adequacy	Samples are ordered via n-gram similarity between	850
797	ensured image caption evaluation. In <i>Proceedings of</i>	the predictions of both models, listing the most	851
798	<i>the IEEE/CVF Conference on Computer Vision and</i>	different predictions first, allowing for better em-	852
799	<i>Pattern Recognition</i> , pages 14050–14059.	phasizing items with evident differences first. This	853
800	Adina Williams, Nikita Nangia, and Samuel Bowman.	is calculated by considering the top 5 beams of	854
801	2018. <a href="#">A broad-coverage challenge corpus for sen-</a>	BLIP as reference texts and the top 5 beams of	855
802	<a href="#">tence understanding through inference.</a> In <i>Proceed-</i>	BLIP+ <i>MOCHA</i> as candidate sentences; we then	856
803	<i>ings of the 2018 Conference of the North American</i>	compute the average BLEU (Papineni et al., 2002)	857
804	<i>Chapter of the Association for Computational Lin-</i>	score between each candidate and all references.	858
805	<i>guistics: Human Language Technologies, Volume 1</i>		
806	<i>(Long Papers)</i> , pages 1112–1122. Association for	<b>B Additional Details</b>	859
807	Computational Linguistics.		
808	Zhenlin Xu, Yi Zhu, Tiffany Deng, Abhay Mittal, Yan-	<b>B.1 MOCHA Implementation Details</b>	860
809	bei Chen, Manchen Wang, Paolo Favaro, Joseph	As discussed in Rennie et al. (Rennie et al., 2017),	861
810	Tighe, and Davide Modolo. 2023. Challenges of	we reduce variance in gradient estimates by shifting	862
811	zero-shot recognition with vision-language mod-	the reward function to have zero mean; we apply	863
812	els: Granularity and correctness. <i>arXiv preprint</i>	this to the reward function before adding the KL	864
813	<i>arXiv:2306.16048.</i>	penalty. To perform this shifting, we subtract the	865
814	Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao	sample mean of this reward (without KL penalty)	866
815	Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun,	among all predictions for a given image in a mini-	867
816	and Enhong Chen. 2023. <a href="#">Woodpecker: Hallucination</a>	batch.	868
817	<a href="#">correction for multimodal large language models.</a>	During each training iteration, we build mini-	869
818	Peter Young, Alice Lai, Micah Hodosh, and Julia Hock-	batches by selecting 10 images and then generat-	870
819	enmaier. 2014. <a href="#">From image descriptions to visual</a>	ing 10 predictions per image (hence 100 image-	871
820	<a href="#">denotations: New similarity metrics for semantic in-</a>	prediction pairs total). We use nucleus sam-	872
821	<a href="#">ference over event descriptions.</a> <i>Transactions of the</i>	pling (Holtzman et al., 2019) with $p = 0.9$ and	873
822	<i>Association for Computational Linguistics</i> , 2:67–78.	temperature $t = 1.2$ , and we cap generations to	874
823	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q		
824	Weinberger, and Yoav Artzi. 2019. Bertscore: Eval-	<sup>1</sup> Using the same pretrained NLI model described in the	
825	uating text generation with bert. <i>arXiv preprint</i>	main paper.	
826	<i>arXiv:1904.09675.</i>		

be at most 40 tokens. We apply PPO reinforcement with clipping parameter  $\epsilon = 0.2$ . For our reward function, we use coefficients  $\alpha = 0.5$  and  $\beta \in [0.004, 0.06]$  (depending on the model optimized).

During *MOCHA* training, we freeze the image encoder of all models, training the text encoder components alone. For BLIP-Large and BLIP-Base we use gradient clipping of 5, learning rate of  $1e-6$  and 4 PPO steps in each iteration. BLIP-2 is trained with low rank adapters (LoRA) over the keys and values of the decoder attention layers (Hu et al., 2021) with a learning rate of  $1e-6$ . GIT-base is trained with a learning rate of  $1e-5$  with 4 PPO steps and gradient clipping of 5.

All model checkpoints are taken from the Hugging Face Model Hub<sup>2</sup>:

- salesforce/blip-image-captioning-large
- salesforce/blip-image-captioning-base
- salesforce/blip2-opt-2.7b-coco
- microsoft/git-base-coco

We train these models for the following number of iterations: 350 for BLIP-B, 1200 for BLIP-L, 3400 for BLIP-2, and 600 for GIT-B.

## B.2 OpenCHAIR Implementation details

**Generating Diverse Captions** We start by parsing all objects in MS-COCO’s human-annotated captions by first identifying nouns via syntactic parsing<sup>3</sup>. We then filter these for highly concrete nouns, by using the values recorded by Hessel et al. (Hessel et al., 2018) with threshold 4.5. We used these objects, coupled with their corresponding captions, to prompt an instruction-tuned LLM<sup>4</sup> to rephrase the captions with different objects. We used stochastic sampling with top-p of 0.9 and temperature of 0.6 for this LLM generation. While this stage increases the object diversity, we notice that the output still includes many common objects that have a significant overlap with those in MS-COCO. To overcome this issue, we filter out all captions that do not include rare objects, defining an object as rare if its appearance frequency in the dataset is in the lowest 10th percentile. The remaining captions are used as few-shot examples for a LLM<sup>5</sup>

<sup>2</sup><https://www.huggingface.co/models>

<sup>3</sup>Using the *en\_core\_web\_md* pipeline from the SpaCy (Honnibal and Montani, 2017) library.

<sup>4</sup>meta-llama/Llama-2-70b-chat-hf (4-bit quant.)

<sup>5</sup>meta-llama/Llama-2-13b

(base, not instruction-tuned) to generate new captions, to further increase diversity. We used 10 few shot example for each generated caption, and text is generated using sampling with temperature 0.8. We generate 2,000 captions from the LLM and feed them as prompts to the text-to-image generation model Stable Diffusion XL (Podell et al., 2023), which generates a single image for each caption. For image generation, we use 40 sampling steps and guidance scale of 10. We also employ negative prompting using the prompt “*unclear, deformed, out of image, disfigured, body out of frame*” to encourage generation of clear objects in the output images.

## Evaluation on the OpenCHAIR Benchmark

Evaluating a captioning model on *OpenCHAIR* is performed as follows: First, all the objects in the caption generated by the captioning model are extracted using the parsing method described in the previous paragraph. For each detected object, an LLM<sup>4</sup> is prompted to determine whether the object is in the GT caption or not using the prompt: “<s>[INST] An image has the following caption: “<input caption>”. Does the image contain the following object? “<input object>”. Answer yes/no/unsure. The answer is: [/INST]”. We use greedy decoding for this stage. Objects for which the LLM answers “no” are counted as hallucinations and objects for which the LLM answers “yes” are counted as existing objects. We ignore objects that receive any other response, and report that the amount of such objects are <2% of the total objects considered. Finally, the *OpenCHAIR* hallucination rate is calculated as  $OCH := n_h / (n_h + n_e)$ , where  $n_h$  is the number of hallucinated objects and  $n_e$  is the number of existing objects.

## B.3 LURE Comparison

To evaluate LURE (Zhou et al., 2024) in our setup, we followed the authors’ instructions<sup>6</sup> and applied their pre-trained model (YiyangAiLab/LURE, over MiniGPT-4 with VICUNA-13b) to our predicted captions. BLIP’s predictions (with beam search decoding, 5 beams) were supplied to LURE’s revisor along with the probabilities of each predicted token for the highest scoring beam. After additional parsing, LURE revised BLIP-L’s sentences, which we then evaluated with various metrics.

An example of this procedure is given below:

<sup>6</sup><https://github.com/YiyangZhou/LURE/blob/main/README.md>

<b>BLIP2</b>	Pred = 'E' Pred = 'H'	
GT = 'E'	332	42
GT = 'H'	0	54
<b>BLIP-L</b>	Pred = 'E' Pred = 'H'	
GT = 'E'	353	44
GT = 'H'	0	31
<b>GIT-B</b>	Pred = 'E' Pred = 'H'	
GT = 'E'	325	36
GT = 'H'	1	66
<b>OFA-L</b>	Pred = 'E' Pred = 'H'	
GT = 'E'	336	45
GT = 'H'	1	46

Table 3: **Human Evaluation of OpenCHAIR Benchmark.** The tables illustrate a correlation measurement between *OpenCHAIR*’s automatic hallucination annotations (Pred) and manual human hallucination annotations (GT). ‘E’, ‘H’ stand for ‘object **E**xists’, ‘object **H**allucinated’, respectively. BLIP2, BLIP-L, GIT-B and OFA-L stand for BLIP2-2.7b, BLIP-Large, GIT-Base, OFA-Large, all fine-tuned for image-captioning over COCO.

<b>BLIP2</b>	Pred = 'E' Pred = 'H'	
GT = 'E'	416	3
GT = 'H'	4	5
<b>BLIP-L</b>	Pred = 'E' Pred = 'H'	
GT = 'E'	413	2
GT = 'H'	4	9
<b>GIT-B</b>	Pred = 'E' Pred = 'H'	
GT = 'E'	412	1
GT = 'H'	3	12
<b>OFA-L</b>	Pred = 'E' Pred = 'H'	
GT = 'E'	418	2
GT = 'H'	3	5

Table 4: **Human Evaluation of CHAIR Benchmark.** The tables illustrates a correlation measurement between *CHAIR*’s automatic hallucination annotations (Pred) and manual human hallucination annotations (GT). ‘E’, ‘H’ stand for ‘object **E**xists’, ‘object **H**allucinated’, respectively. BLIP2, BLIP-L, GIT-B and OFA-L stand for BLIP2-2.7b, BLIP-Large, GIT-Base, OFA-Large, all fine-tuned for image-captioning over COCO.

- Original caption from BLIP, masked by LURE’s algorithm: *a woman standing in a kitchen with [IDK]*
- LURE instruction: *According to the picture, remove the information that does not exist in the following description:*
- LURE output revision: *In the image, a smiling woman is standing in a kitchen with various food items spread out on a table. She seems to be posing and enjoying the scene, perhaps taken during a cooking session. The table is filled with food, including several pieces of sushi placed near the center, along with some pasta on the left side, cucumbers near the right side, and a couple of apples towards the back.*

#### B.4 TLC-A Comparison

In order to compare our method to TLC-A (Petryk et al., 2023), we received code from its authors and implemented it in our setup. TLC-A is a decoding-time method applied to auto-regressive captioning models, and in our setting we apply it to different models (e.g. BLIP-Large) than those tested by Petryk et al (e.g. OFA). Of particular note is that

TLC-A requires selecting a threshold confidence value, which is used in the decoding phase to re-rank generated beams according to the confidence assigned to COCO object tokens. Petryk et al. recommend calibrating this threshold using the COCO validation set to achieve a precision level of at least 99%; however, in our experiments we find that this value cannot be achieved by the models we consider without sacrificing most of the recall, as illustrated in Figure 11. Therefore, we instead use the COCO validation set to select the best-performing threshold with respect to the CHAIR metric, as shown in Table 5. The selected confidence threshold is 0.33 and it achieves a precision of 98.3% and a recall of 84% over the validation set.

## C Additional Results

### C.1 Full Quantitative Results

We show in Table 6 the full results, comparing the *MOCHA* optimized models (marked by +M) to the baselines (Figure 7 was prepared using this data).

<sup>6</sup>Reference ground truth captions: *Painting of oranges, a bowl, candle, and a pitcher* (left) and *A giraffe grazing on a tree in the wilderness with other wildlife* (right).



TH	P	R	B@4↑	C↑	CH <sub>i</sub> ↓	CH <sub>s</sub> ↓	$\bar{p}$ ↓	BSc↑
-	-	-	41.5	138.4	2.3	3.5	0.246	0.679
0.10	0.978	0.99	41.4	138.0	2.2	3.38	0.246	0.677
0.21	0.980	0.94	41.4	137.7	2.1	3.14	0.243	0.677
<b>0.33</b>	<b>0.983</b>	<b>0.84</b>	<b>41.2</b>	<b>137.5</b>	<b>1.91</b>	<b>2.82</b>	<b>0.243</b>	<b>0.676</b>
0.52	0.986	0.61	41.1	136.7	1.97	2.9	0.242	0.675
0.56	0.988	0.55	41.2	136.8	1.94	2.86	0.243	0.675
0.94	1	0.01	41.4	137.7	2.21	3.32	0.247	0.677

Table 5: **Selecting a threshold for TLC-A.** We evaluate TLC-A with different thresholds (as described by Petryk et al. (Petryk et al., 2023)) over the COCO caption Karpathy validation set. In the first row we have BLIP without TLC-A. We indicate the selected threshold which achieves the best CHAIR scores overall in **bold**. B@4, C, CH<sub>i</sub>, CH<sub>s</sub>, BSc,  $\bar{p}$  denote BLEU-4, CIDEr, CHAIR instance and CHAIR sentence, BERTScore, and NLI  $p(\text{contr.})$  metrics respectively. P, R are the precision and recall that each threshold (for predicted object confidences) achieves over the validation set.

Model	B@4↑	C↑	CH <sub>i</sub> ↓	CH <sub>s</sub> ↓	OCH↓	$\bar{p}$ ↓	BSc↑
BLIP-B	24.8	87.5	2.6	2.8	13.4	0.206	0.557
BLIP-B+M (ours)	<b>26.0</b>	<b>91.3</b>	<b>2.2</b>	<b>2.5</b>	<b>11.4</b>	<b>0.176</b>	<b>0.576</b>
BLIP-L	41.5	138.4	2.3	3.5	16.4	0.244	0.679
BLIP-L + TLC-A	41.3	137.4	<b>1.9</b>	<b>2.8</b>	16.4	0.241	0.676
BLIP-L+M (ours)	<b>41.9</b>	<b>139.6</b>	2.1	3.1	<b>15.4</b>	<b>0.206</b>	<b>0.682</b>
BLIP2	43.4	<b>144.3</b>	1.7	2.6	14.7	0.207	<b>0.684</b>
BLIP2+M (ours)	<b>44.0</b>	<b>144.3</b>	<b>1.4</b>	<b>2.3</b>	<b>14.5</b>	<b>0.199</b>	<b>0.684</b>
GIT-B	38.7	128.1	4.2	2.9	21.3	0.284	0.656
GIT-B+M (ours)	<b>39.0</b>	<b>128.4</b>	<b>3.9</b>	<b>2.7</b>	<b>19.6</b>	<b>0.221</b>	<b>0.657</b>

Table 6: **Quantitative results** for state-of-the-art VLM models on the COCO Caption Karpathy test set. +M refers to *MOCHA*. BSc and  $\bar{p}$  denote BERTScore and NLI contradiction probability rewards. B@4, C, CH, OCH, BSc and  $\bar{p}$  denote BLEU-4, CIDEr, CHAIR (i for instance, s for sentence), OpenCHAIR, BERTScore, and NLI  $p(\text{contr.})$  metrics respectively. All results are generated by using their officially provided checkpoints and hyperparameters. Best results are shown in **bold**.

## C.2 Comparisons of *OpenCHAIR* and CHAIR

In Tables 3–4 we provide full numeric results for our human evaluation of *OpenCHAIR* and CHAIR across a variety of captioning model predictions, as we discuss in the main paper.

In Figure 14, we illustrate the number of unique object types found in these benchmarks. We note that *OpenCHAIR* contains a much larger diversity of object types, even when considering the full contents of CHAIR’s synonym list.

## C.3 Additional Ablations

**Reward Ablations.** In Table 9, we provide numeric results for ablating the fidelity and adequacy

terms in our reward function. As discussed in the main paper, removing either of these reward terms leads to a degradation with respect to either hallucinations or textual quality, while using both together displays a synergistic effect with hallucinations reduced (as reflected by metrics such as CHAIR) while preserving or even improving caption quality (as reflected by general textual quality metrics such as BLEU-4). We also show a qualitative illustration of these results in Figure 15.

We demonstrate the effect of our KL penalty in the reward function by performing *MOCHA* optimization without this term. As can be observed in the fifth row of Table 7, optimization without this penalty improves the NLI-based reward  $\bar{p}$  while degrading other measures of text quality (including non-optimized metrics like CIDEr). We hypothe-

<sup>2</sup>Reference ground truth captions: *A car with some surfboards in a field* (left) and *A boy holding umbrella while standing next to livestock* (right).

Model	OCH ↓	B@4 ↑	C ↑	CH <sub>i</sub> ↓	CH <sub>s</sub> ↓	$\bar{p}$ ↓	BSc ↑
BLIP-L	0.270	41.5	138.4	2.3	3.5	0.244	0.679
BLIP-L+M	0.259	41.9	139.6	2.1	3.1	0.206	0.682
$-r_f$	0.267	43.0	142.3	2.8	4.4	0.249	0.691
$-r_a$	0.257	41.1	132.9	1.5	2.3	0.174	0.66
$-r_{kl}$	0.241	27.6	98.9	1.4	1.9	0.135	0.62
$-ppo$	0.287	39.4	127.6	2.5	3.76	0.212	0.664

Table 7: **Additional ablation results.** We ablate the effect of the KL penalty reward  $r_{kl}$  and the selection of PPO algorithm. As seen above, removing  $r_{kl}$  causes the model to over-optimize the fidelity reward ( $\bar{p}$ ), while replacing PPO with the simpler SCST algorithm (described in Section C.3) leads to instabilities that degrade performance across metrics.

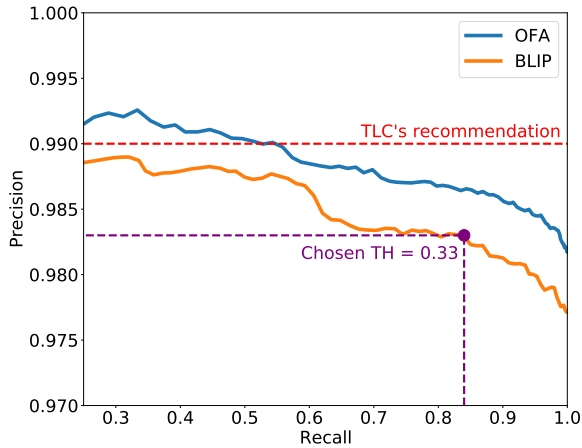


Figure 11: **Precision-recall curve for selecting TLC-A threshold.** As detailed in (Petryk et al., 2023), we compute a precision-recall curve over the predicted object confidences. As illustrated above, the 99% precision threshold recommended by Petryk et al. (Petryk et al., 2023) cannot be achieved by BLIP-Large on the COCO Karpathy validation set. Hence, in our setting we must adjust the threshold to find a reasonable balance between precision and recall.

$\emptyset$	<i>a painting of oranges and a silver pitcher on a table</i>	<i>two giraffes eating leaves from a tree</i>
$-r_{kl}$	<i>a painting of some items</i>	<i>some giraffes in the field</i>
$r$	<i>a painting of a pitcher, oranges, and a candle on a table</i>	<i>a giraffe eating leaves from a tree in a field</i>

Figure 12: **Ablating the KL-penalty reward.** Above we show captions sampled from various models: the initial model (BLIP-Large) before optimization ( $\emptyset$ ), the model with *MOCHA* optimization applied and KL penalty ablated ( $-r_{kl}$ ), and an optimized model with our full reward function ( $r$ ). As is seen above, while the base model outputs various hallucinations (e.g. *a silver pitcher*), the model optimized without KL penalty outputs generic texts without adequate detail, due to over-optimization of the fidelity objective. Optimizing with the full reward function yields captions that are both descriptive and consistent with the input condition.

	LLaVa-RLHF	BLIP-L+ <i>MOCHA</i>
	<i>A man sitting on a chair with a large stuffed animal, specifically a teddy bear, on his lap</i>	<i>a man sitting on a chair holding a large stuffed animal</i>

Figure 13: **LLaVa-RLHF vs. *MOCHA*.** We illustrate that RLHF training does not necessarily solve the hallucination problem of VLM models by showing a generation produced by LLaVa-RLHF (Sun et al., 2023a) compared to BLIP+*MOCHA*. For both models, we use the prompt “a photography of” for generation. See Table 10 for a quantitative comparison.

Model	B@4↑	M↑	C↑	CH <sub>s</sub> ↓	CH <sub>i</sub> ↓
<b>Dedicated</b>					
UD-L+Occ <sub>XE</sub>	33.9	27.0	110.7	5.9	3.8
UD-L+Occ <sub>SC</sub>	37.7	28.7	125.2	5.8	3.7
CIIC <sub>XE</sub>	37.3	28.5	119.0	5.3	3.6
CIIC <sub>SC</sub>	40.2	29.5	133.1	7.7	4.5
COSNet <sub>XE</sub>	39.1	29.7	127.4	4.7	3.2
COSNet <sub>SC</sub>	<u>42.0</u>	30.6	<u>141.1</u>	6.8	4.2
<b>End-to-end</b>					
BLIP	41.5	<u>31.1</u>	138.4	<u>3.5</u>	<u>2.3</u>
BLIP-2	<b>43.4</b>	<b>31.7</b>	<b>144.3</b>	<b>2.6</b>	<b>1.7</b>

Table 8: **Older dedicated methods for reduced-hallucination captioning vs. end-to-end modern VLMs for image captioning.** Results are given on the Karpathy test split of MS-COCO dataset, including closed-vocabulary hallucination metrics as commonly reported by such dedicated methods. B@4, C, M, CH denote BLEU-4, CIDEr, METEOR, and CHAIR metrics respectively. We see that older, dedicated methods with weaker backbones are outperformed by modern VLMs on all metrics, including the smaller BLIP(-Large) and the larger BLIP-2(-2.7B). XE and SC indicate cross-entropy and SCST (RL) optimization respectively. Best and second-best metric values are shown in **bold** and underlined text respectively.

size that allowing the model to freely deviate from its initial distribution encourages it towards a degenerate solution with respect to  $\bar{p}$ , which may be the easiest reward term to over-optimize in an unconstrained setting. This is also reflected qualitatively as seen in Figure 12. As illustrated in the figure, captions generated by the model trained without the KL penalty ( $-r_{kl}$ ) do not contradict the image, but rather contain generic text (e.g. *a painting with some items*), lacking adequate detail. By contrast, optimizing with the KL penalty reward yields captions that are both descriptive and consistent with the input condition, reflected in the improved scores across metrics in Table 7 and the quality of predictions of the full reward model ( $r$ ) in Figure 12. This is attributed to the ability of the KL penalty to mitigate over-optimization, which benefits both optimized rewards.

**PPO Ablation.** We also ablated the selection of RL algorithm, by replacing PPO with the SCST algorithm upon which it is based (noting that SCST is the common name for the REINFORCE algorithm in the context of image captioning) (Sutton and Barto, 2018; Schulman et al., 2017; Rennie et al., 2017). As is seen in Table 7, PPO outper-

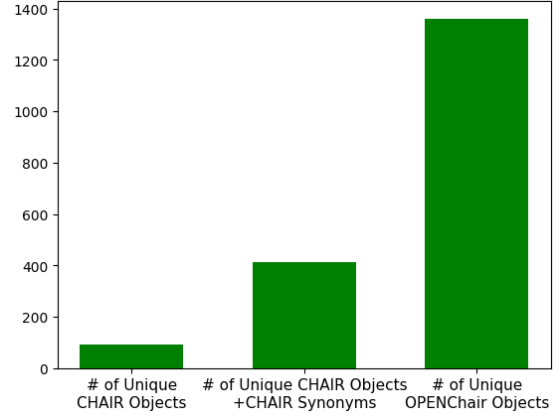


Figure 14: **Object Type Coverage, CHAIR vs. OpenCHAIR.** We display the object type coverage of CHAIR (over MS-COCO) and *OpenCHAIR*, measured as the number of unique objects. In OPENChair, objects are found using the parsing method described in Section B.2. As can be observed, the proposed benchmark has significantly greater coverage of different objects.

Model	B@4↑	C↑	CH <sub>i</sub> ↓	CH <sub>s</sub> ↓	$\bar{p}$ ↓	BSc ↑
BLIP	41.5	138.4	2.3	3.5	0.246	0.679
BLIP+M	41.9	139.6	2.1	3.1	0.206	0.682
$-r_f$	43.0	142.3	2.8	4.4	0.249	0.691
$-r_a$	41.1	132.9	1.5	2.3	0.174	0.66

Table 9: **Reward Ablation.** We ablate the effect of the fidelity  $r_f$  and adequacy  $r_a$  terms in our reward function, finding that using each alone significantly degrades performance with respect to hallucinations or textual quality.

forms SCST across metrics, consistent with prior work on PPO finding that it avoids instabilities during optimization that may allow it to converge to a more optimal solution (Schulman et al., 2017; Ouyang et al., 2022; Ziegler et al., 2020).

#### C.4 Additional Comparisons

**Comparison to Dedicated Models** In Table 8 we provide full numeric results for older dedicated models compared to a modern VLM without further optimization, showing that they are outperformed by all metrics.

**Comparison to RLHF-Tuned VLMs.** LLaVa-RLHF (Sun et al., 2023a) is a concurrent work, which aims to reduce hallucinations in instruction tuned models using factually-grounded RLHF. In Table 10, we provide a quantitative comparison between LLaVa-RLHF and BLIP+*MOCHA* over 100 samples of the OPENChair benchmark. For LLaVa-RLHF decoding we use both stochas-






		
$\emptyset$	<i>This is a picture of a large old fashioned car that was parked by a group of people</i>	<i>People at festival standing around in open field</i>
$-r_f$	<i>A car parked in the grass with a surfer standing near it</i>	<i>A woman standing next to a herd of animals with an umbrella</i>
$-r_a$	<i>Spectators could enjoy the old fashions of the fifties</i>	<i>That are some very nice people who are very fun to view them</i>
$r$	<i>A vintage car parked on a field next to people</i>	<i>A young man with a large umbrella next to a herd of animals</i>

Figure 15: **Ablating our multi-objective reward function.** Above we show captions sampled from models with different reward functions. Top row depicts the initial model (before optimization). As can be seen in the table, generations of the base model ( $\emptyset$ ) and the model trained without the fidelity objective ( $-r_f$ ) contain various hallucinations that contradict the image, like stating that the car was *parked by a group of people*, confusing between an ordinary person and a *surfer*, and stating that the boy is a *woman*. In contrast, those from the model without the adequacy objective ( $-r_a$ ) are generic and neutral with respect to the image (without explicitly contradicting it), e.g. the abstract statement about the *spectators enjoying the old fashions of the fifties*. At last, optimizing for both ( $r$ ) yields captions that are both descriptive and consistent with the input condition, similar to the reference captions<sup>2</sup> that were provided by human annotators.

tic sampling with the default parameters recommended by the authors, as well as greedy sampling (as beam search is not implemented for LLaVa-RLHF). For a fair comparison, we use greedy decoding for BLIP+MOCHA as well. As LLaVa-RLHF tends to generate long paragraphs which follow an image description with subjective commentary, we terminate generation after a single sentence, which usually corresponds to an image caption. The instruction given to LLaVa-RLHF is “describe the image briefly”. As seen in the table, our method outperforms LLaVa-RLHF by this measure of open-vocabulary hallucinations. This is further seen in Figure 13, which shows example captioning predictions for these models, illustrating that LLaVa-RLHF may be more prone to hallucinations.

**Evaluation over Flickr30K dataset.** We per-

Model	OCH ↓
LLaVa-RLHF <sub>S</sub>	0.396
LLaVa-RLHF <sub>G</sub>	0.401
BLIP-L+M <sub>G</sub>	<b>0.360</b>

Table 10: OPENChair comparison between LLaVa-RLHF and BLIP-L+MOCHA over 100 random samples. For LLaVa-RLHF, S stands for stochastic sampling with default parameters, and G stands for greedy decoding (as beam search is not implemented for LLaVa-RLHF). For fair comparison, we also apply greedy decoding to BLIP-L+MOCHA.

Model	B@4↑	C↑	$\bar{p}$ ↓	BSc ↑
BLIP	<b>29.0</b>	73.2	0.335	0.603
BLIP+M	28.9	<b>73.6</b>	<b>0.296</b>	<b>0.607</b>

Table 11: **Evaluation over Flickr30K dataset.** We perform a zero-shot evaluation of BLIP-Large with and without MOCHA (performed on COCO) on an additional dataset. As seen above, improvements to the optimized metrics ( $\bar{p}$  and BERTScore) transfer to the new dataset, while other text quality metrics have similar values before and after MOCHA-tuning, suggesting that overall text quality is generally preserved.

form a zero-shot generalization test by evaluating a MOCHA-tuned model on an additional dataset (different from COCO upon which the model was MOCHA-tuned). In Table 11 we can see that the model with MOCHA fine-tuning shows an improvement in metrics (NLI and BERTScore) that were optimized on the training data from COCO. Furthermore, we see that non-optimized text quality metrics have similar values between both models, suggesting that MOCHA tuning generally preserves overall text quality. Supporting this quantitative evaluation, we provide detailed qualitative results on the Flickr30K dataset in the attached visualization tool.