ENHANCING FACTUALITY IN DETAILED IMAGE CAP TIONING WITH LLM-MLLM COLLABORATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Multimodal large language models (MLLMs) capable of interpreting images can generate highly detailed and extensive captions, owing to their advanced language modeling capabilities. However, the captions they produce frequently contain hallucinations. Furthermore, our empirical analysis reveals that existing hallucination detection methods are less effective in detailed image captioning tasks. We attribute this to the increasing reliance of MLLMs on their own generated text, rather than the input image, as the sequence length grows. To address this issue, we propose a novel corrector-based method that decomposes a given caption into atomic propositions, evaluates the factuality of each unit, and revises the caption accordingly. Our method is training-free and can be applied in a plug-and-play manner to any captioning model. Additionally, we introduce an evaluation framework and a benchmark dataset to facilitate the systematic analysis of detailed captions. Our experiments demonstrate that existing approaches to improve the factuality of MLLM outputs may fall short in detailed image captioning tasks. In contrast, our proposed method significantly enhances the factual accuracy of captions, even improving those generated by GPT-4V. Finally, we highlight a limitation of VQAcentric benchmarking by demonstrating that an MLLM's performance on VQA benchmarks may not correlate with its ability to generate detailed image captions.

027 028 029

025

026

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Numerous image captioning methods utilizing deep neural networks (DNNs) have been proposed
(Vinyals et al., 2015; Xu et al., 2015). However, they are generally limited to generating short and
concise captions, which constrains their broader application in real-world scenarios. For instance, in
cases such as assistance for visually impaired individuals, where it is necessary to provide highly
detailed descriptions of the scene in front of the user, these methods may not be suitable.

Following the recent success of large language models (LLMs) (Brown et al., 2020), there have been 037 attempts to use not only text but also information from other modalities as input to LLMs. Notably, 038 many studies have explored multimodal large language models (MLLMs) that incorporate visual information (Li et al., 2023a; Dai et al., 2023; Liu et al., 2024b). These models have demonstrated significantly superior performance compared to traditional models in tasks such as visual question 040 answering (VQA) and captioning (Liu et al., 2024a). In particular, MLLMs, leveraging the advanced 041 language capabilities of LLMs, are able to generate much longer and more detailed captions than 042 conventional captioning models. However, these generated captions frequently contain inaccurate 043 information, including descriptions of objects that are not present in the input image (Leng et al., 044 2024). Such hallucination problems hinder the practical application of MLLMs in real-world settings.

Three major approaches have been recently proposed to improve the factuality of MLLM outputs: (i) Decoding-based methods (Leng et al., 2024) reduce the probabilities of hallucination-related tokens during the model's decoding process without requiring additional training; (ii) Training-based methods (Liu et al., 2023a) further train the models on curated multimodal datasets to ensure they generate only accurate responses; (iii) Corrector-based methods (Zhou et al., 2024) employ a corrector model that detects and either removes or revises hallucinations present in the model's responses.

In this paper, we propose a novel corrector-based method called Visual Factuality EnhanceR
 (V-FactER). Unlike existing approaches that require training a corrector (Lee et al., 2024), V-FactER improves the factuality of detailed image captions by leveraging the collaboration between an LLM

and MLLM, without the need for additional training. Moreover, unlike methods that target specific
types of hallucinations (Li et al., 2023b; Zhou et al., 2024), our approach does not pre-define the
hallucination types, allowing it to address a broader range of issues. The method proceeds as follows:
(i) an LLM decomposes a given detailed caption into atomic propositions; (ii) an MLLM verifies the
truthfulness of each atomic proposition based on the corresponding image; and (iii) the LLM revises
the caption accordingly. Our design is particularly motivated by the observation that, as the length of
a model's response increases, hallucinations generated later in the sequence become more difficult
for existing methods (Wang et al., 2023; Zhou et al., 2024) to detect.

062 Evaluating the factuality of detailed captions is not straightforward. Through experiments, we 063 demonstrate that conventional caption evaluation metrics such as BLEU (Papineni et al., 2002), 064 ROUGE (Lin, 2004), METEOR (Banerjee & Lavie, 2005), and CIDEr (Vedantam et al., 2015), as well as recently proposed methods (Hessel et al., 2021; Petryk et al., 2024), fail to accurately assess 065 the factuality of detailed captions. To address this issue, we propose a novel GPT-based method 066 for factuality evaluation and validate its effectiveness through experiments. Even if a caption 067 contains factual information, however, it may still be considered inadequate if it does not sufficiently 068 capture the visual information. To measure the coverage of captions, we construct a detailed VQA 069 dataset through a collaboration between humans and an AI agent (Achiam et al., 2023). If a caption fully encapsulates the information of a given image, questions about the image should be 071 answerable accurately using only the caption, without referencing the image itself. 072

Our experiments surprisingly reveal that methods designed to improve the factuality of MLLMs, which have proven effective in tasks like VQA (Huang et al., 2024), may be ineffective for detailed image captioning tasks that require longer responses. In contrast, V-FactER significantly enhances the factuality of captions and can be applied in a plug-and-play manner to any captioning model. Our experiments further demonstrate that this improvement extends to captions generated by the state-of-the-art closed model, GPT-4V (Achiam et al., 2023). Finally, we highlight an issue with the current VQA-centric benchmarking (Duan et al., 2024) by showing that an MLLM's performance on VQA benchmarks may not correlate with its ability to generate detailed image captions.

In summary, our **contributions** are as follows:

- We demonstrate that existing hallucination detection methods may perform worse as MLLM response length increases, and we propose a method that can circumvent this issue.
- We introduce V-FactER, a method that significantly enhances the factuality of given detailed image captions. V-FactER is a pipeline that leverages a pre-trained LLM and MLLM.
- We propose an evaluation framework and benchmark dataset that overcome the limitations of existing caption evaluation methods and enable the systematic analysis of detailed image captions.
- We show that while existing methods designed to improve the factuality of MLLM responses may be ineffective for detailed image captioning tasks, V-FactER significantly improves their factuality.
- Our experiments demonstrate that current VQA benchmarks fail to reliably capture the potential of MLLMs in real-world applications, such as visual assistants for the visually impaired.
- 092

083

084

085

087

088

090

091

093 094

2 RELATED WORK

Multimodal large language models. LLMs that process inputs from multiple modalities, including 095 text and other types of data, are referred to as multimodal LLMs (Yin et al., 2023a). Among these, 096 LLMs that handle visual input have been the most actively researched, and the MLLMs discussed 097 in this paper are focused on this category. Research on these models primarily explores methods 098 for fusing the output of an independent vision encoder into the input of an LLM. The BLIP models (Li et al., 2023a; Dai et al., 2023) align the frozen vision encoder and LLM using a lightweight 100 transformer (Vaswani, 2017) called Q-Former. The trainable input tokens of the Q-Former interact 101 with the output tokens from the vision encoder through cross-attention, transforming them into 102 input tokens for the LLM. The LLaVA models (Liu et al., 2024b;a) use a simple MLP connector to 103 align the vision encoder with the LLM. All output tokens from the vision encoder, passed through 104 the MLP connector, are used as input to the LLM. The vision encoder's parameters remain fixed 105 during the training of the MLP connector and the LLM. Unlike existing MLLMs, the InternVL models (Chen et al., 2024c;b) have demonstrated the effectiveness of increasing the size of both the 106 vision encoder and the vision-language connector. They utilize a 6-billion parameter vision encoder 107 and an 8-billion parameter vision-language connector. The connector is obtained by fine-tuning the



112

113

114 115

116

121

122

124



Figure 1: The process of generating a data sample for evaluating the performance of hallucination detection methods in detailed image captioning tasks. Human annotators identify and label object hallucinations within the caption generated by LLaVA-NeXT (Liu et al., 2024a) for an input image.

pre-trained multilingual LLaMA (Cui et al., 2023). Despite the many advancements in open-source MLLMs, closed-source MLLMs such as GPT-4V or GPT-40¹ still outperform them significantly. As a result, these GPT models represent the upper bound performance in benchmarks and are commonly 123 used to evaluate MLLMs (Petryk et al., 2024). In our work, we demonstrate that captions generated by GPT-4V can be improved using our method, and we use GPT-40 to assess the factuality of captions. 125

126 MLLM hallucinations and mitigation strategies. MLLMs sometimes generate inaccurate re-127 sponses. For example, they may incorrectly describe the characteristics of objects in an input image, 128 misrepresent relationships between objects, or even describe objects that do not exist. To mitigate 129 these hallucination problems, decoding-based methods identify factors that induce hallucinations and 130 apply penalties to the probabilities of tokens that are likely to be hallucinations during the decoding 131 process. For instance, VCD (Leng et al., 2024) induces hallucinations using corrupted images, while 132 OPERA (Huang et al., 2024) leverages the correlation between high attention weights assigned to a 133 few summary tokens and hallucinations. Training-based methods focus on exploring training data that can suppress the generation of hallucinations. Liu et al. (2023a) demonstrated that hallucinations 134 can be alleviated by incorporating negative samples-descriptions that explicitly state the absence of 135 certain objects in a given image-into visual instruction tuning datasets. Corrector-based methods 136 (Zhou et al., 2024; Lee et al., 2024) detect, remove, and revise hallucinations present in MLLM 137 responses by using a corrector model. This model is obtained by supervised fine-tuning a pre-trained 138 MLLM. The corrector model then revises the initial response based on the given image. 139

140

Caption evaluation methods. Since short image captions are relatively easy to obtain reference 141 captions for, we can use matching-based caption evaluation methods (Hossain et al., 2019) to assess 142 them. However, for long and detailed captions generated by MLLMs, the number of reference 143 captions required for such evaluations becomes exceedingly large. Thus, it becomes impractical to 144 evaluate detailed captioning using traditional approaches. Hessel et al. (2021) proposed CLIPScore, a 145 reference-free evaluation method. CLIPScore measures the distance between an image and its caption 146 within the pre-trained joint representation space of CLIP (Radford et al., 2021). Additionally, the 147 authors introduced RefCLIPScore, which uses both the image and reference captions within that same representation space. Chan et al. (2023) addressed the limitations of matching-based methods 148 by utilizing an LLM. The LLM-based metric they proposed, CLAIR, assigns scores to captions based 149 on reference captions using an LLM. Similarly, ALOHa (Petryk et al., 2024) detects hallucinations 150 by comparing a generated caption and its reference information through the use of an LLM. 151

152 153

3 METHOD

155 In this paper, we propose a new corrector-based method. Corrector-based methods typically detect and 156 remove or revise hallucinations within model responses. Unlike existing approaches, which obtain the 157 corrector model through training, our method employs collaboration between a pre-trained MLLM 158 and LLM. Moreover, in contrast to previous methods that are limited to correcting specific types of 159 hallucinations (Zhou et al., 2024), our approach is free from such constraints. We also propose a dataset and framework for evaluating the detailed image captioning capabilities of an MLLM. Unlike 160

161

https://openai.com/index/hello-gpt-4o/



Figure 2: The hallucination scores of the Confidence and Consistency methods based on object positions within detailed captions. Object hallucinations near the end of the captions (192+) are undetectable by both methods.

existing methods, our proposed evaluation approach allows for the assessment of image captioning models in terms of both factuality and coverage, evaluating each of these aspects separately.

3.1 MOTIVATING OBSERVATIONS

173

174 175 176

177

178 179

181

Here, we examine the performance of existing hallucination detection methods on tasks that require
generating long responses. To facilitate these analyses, we construct a dataset as follows: (i) We
prompt an MLLM with "Describe the given image in a very detailed manner." and collect the model's
responses for a specified image set; (ii) For the convenience of our analysis, we use an LLM to identify
objects that may be hallucinations; (iii) Human annotators then label each parsed object as either a
hallucination or not, based on the corresponding image. We use LLaVA-NeXT (Liu et al., 2024a)
and GPT-4 as the MLLM and LLM, respectively. Figure 1 illustrates the process of constructing
the dataset. To build the dataset, we use a subset of IIW-400 (Garg et al., 2024). We detect object
hallucinations using two of the most widely adopted hallucination detection methods:

- 190 1. **Confidence** (Zhang et al., 2023; Zhou et al., 2024): This method detects hallucinations using the probability p_{obj} predicted for the object token when LLaVA-NeXT generates the caption. For multi-token objects, the product of the token probabilities is used. The hallucination score is defined as $H_{obj} = -\log p_{obj}$. A higher H_{obj} indicates a greater likelihood of hallucination.
- 2. Consistency (Wang et al., 2023; Zhao et al., 2024): This method assumes that hallucinations are more influenced by randomness during decoding. Using stochastic decoding, we have LLaVA-NeXT generate 40 detailed captions per image and count the occurrence t_{obj} of each object in the dataset of Figure 1. The hallucination score is defined as H_{obj} = -log t_{obj}/40.

Figure 2 presents the hallucination scores of each method by the position of objects appearing within 199 detailed captions. The horizontal axis of the graphs represents bins of object token indices, with 200 larger token indices indicating positions closer to the end of the caption. The vertical axis represents 201 the mean and standard deviation of the hallucination scores within each bin. Note that Figure 2a 202 reflects the positions and hallucination scores during greedy decoding, while Figure 2b is derived 203 from the average positions and hallucination scores across 40 stochastic decoding iterations. Figure 2 204 demonstrates that hallucinations generated after the 192nd token are undetectable by the Confidence 205 and Consistency methods. Based on these results, we can infer that existing hallucination detection 206 methods may be ineffective in detecting hallucinations in long detailed captions.

207 Our hypothesis regarding these results is that as 208 MLLM outputs become longer, they become more 209 strongly grounded in the text they generate rather 210 than the given image. In fact, our hypothesis is sup-211 ported by several recent studies. For example, Liu 212 et al. (2024c) demonstrated that as MLLM responses 213 lengthen, the attention weights assigned to image tokens decrease, and Zhong et al. (2024) showed 214 that MLLM responses are significantly influenced by 215 prior dialogue. Based on this hypothesis, we test a

Table 1: Performance comparison	of hallucination
detection methods for the dataset	of Figure 1.

Method	AUROC↑	FPR95↓
Confidence	57.5	95.1
Consistency	73.5	75.6
Object Detector	r 61.5	95.7
Isolation	81.4	71.7



Figure 3: Overview of V-FactER. The decomposer LLM breaks down an initial caption into atomic units. These units are converted into True/False questions and fed into the MLLM along with the image, where each unit is assigned a hallucination score according to Equation (1). Each unit is classified as True or False based on the threshold π , and the corrector LLM then revises the initial caption based on these results.

239 method for determining whether each object is a hallucination by disconnecting it from its con-240 text (Isolation). The Isolation method involves querying the LLaVA-NeXT model with parsed 241 objects using the prompt template, "Is there a {} in the photo?" along with the image. When the 242 probability of the "Yes" token for the object query is $p_{Yes|obj}$, the hallucination score is defined as 243 $H_{\rm obj} = -\log p_{\rm Yes|obj}$. We compare the object hallucination detection performance of the Isolation 244 method with that of the Confidence method, the Consistency method, and a method based on an 245 object detector (**Object Detector**) introduced in recent studies (Yin et al., 2023c; Ge et al., 2024). 246 We measure their detection performance on the dataset of Figure 1 using Area Under the Receiver 247 Operating Characteristic (AUROC) and False Positive Rate at 95% true positive rate (FPR95). Table 1 248 demonstrates that the Isolation method outperforms the others. This suggests that breaking a sentence into smaller units and examining each individually can help detect hallucinations in detailed captions. 249

250

262

235

236 237 238

3.2 VISUAL FACTUALITY ENHANCER

Our motivational observations demonstrate that asking about the presence of objects using a prompt 253 template effectively detects object hallucinations in detailed captions. However, this approach has 254 limitations, as it fails to detect various types of hallucinations. To overcome this limitation, we first decompose each detailed caption into atomic propositions using an LLM. An atomic proposition is a 256 claim or statement that must either be true or false. For example, the caption "A house has a red roof 257 and a chimney" is broken down into "A house has a red roof" and "A house has a chimney." We use 258 an LLM to perform this process, but we allow flexibility in cases where the results do not strictly 259 conform to the definition of an atomic proposition. We then investigate the truth of each decomposed 260 unit using an MLLM. Each unit is converted into a True/False question and independently fed to the 261 MLLM. The hallucination score H(u) for the unit u is defined as follows:

$$H(u) = -\log\left(\min\left(p\left(\text{"True"}|x, Q(u)\right) - p\left(\text{"False"}|x, Q(u)\right), \epsilon\right)\right)$$
(1)

264 p ("True") and p ("False") represent the MLLM's token probabilities for the "True" and "False" 265 tokens, respectively. x and ϵ denote the input image and a very small constant near zero. $Q(\cdot)$ is a 266 function that converts the input text into a True/False question, which we implement by prepending 267 "True or False?" to the input. Each unit is included in either the True set \mathcal{T} or the False set \mathcal{F} , based on 268 its hallucination score. To achieve this, we introduce a hyperparameter π , such that $\mathcal{T} = \{u|H(u) \le \pi\}$ and $\mathcal{F} = \{u|H(u) > \pi\}$. Finally, the initial caption, along with the corresponding sets \mathcal{T} and \mathcal{F} , is provided to an LLM, which corrects the initial caption to ensure it contains only factual information. Under review as a conference paper at ICLR 2025

21.8

13.7

272										
273 Evaluation Metric										
274	Caption	BLEU	ROUGE	METEOR	CIDEr	CLIP-S	RefCLIP-S	CLAIR	ALOHa	Ours
275	Clean	4.2	22.0	13.7	6.4	81.3	75.5	86.9	36.2	62.8
276	Object	4.9	22.3	14.5	4.8	81.0	75.3	85.2	31.5	52.3
277	Attribution	4.1	21.8	13.6	6.2	80.9	75.2	80.0	34.3	60.9
611					< -	01.1	<			

6.7

81.4

75.6

83.5

36.9

51.9

Table 2: Meta-evaluation results across various caption evaluation methods. DOCCI and its synthetic hallucinatory
 captions are used for the meta-evaluation. The highest-rated caption for each method is highlighted in **bold**.

We name this method, which improves the factuality of detailed image captions through the collaboration of a pre-trained LLM and MLLM, **Visual Factuality EnhanceR (V-FactER)**. V-FactER is training-free and can be applied in a plug-and-play manner to any captioning model. Unlike existing methods that can only address predefined types of hallucinations, V-FactER can detect and correct all hallucinations at the atomic unit level. The pipeline of V-FactER is illustrated in Figure 3.

3.3 EVALUATION METHODS

4.1

Relation

278 279

281

282

283

284

285 286

287 288

289

290

291

292

293

294

295

296 297

298

Traditional caption evaluation methods rely on word matching between a predicted caption and its reference captions. This approach works because conventional captioning models generate short captions. However, modern MLLMs produce much longer and more detailed captions, making it impractical to obtain sufficient reference captions for accurate evaluation. Given the enriched content of these image captions, rather than simply evaluating them as good or bad, we aim to assess them systematically by considering two key perspectives:

• Factuality: The degree to which the content of the caption is factual and free from hallucinations.

• Coverage: The extent to which the caption captures the information contained in the image.

We propose evaluation methods for detailed image captions from these two perspectives.

299 Factuality. If a human were to measure the factuality of a text, it would be natural to decompose 300 the text into units that can be classified as true or false, and then calculate the proportion of true 301 units (Maynez et al., 2020). We adopt this approach to measure the factuality of captions, utilizing 302 the state-of-the-art model GPT-40. In our framework, GPT-40 decomposes each caption into atomic propositions and determines their truthfulness based on the corresponding image and reference 303 caption. If the number of atomic propositions judged as true and false are T and F, respectively, the 304 factuality of the caption is defined as $\frac{T}{T+F}$. This approach enables reliable factuality evaluation using 305 only a single reference caption, unlike conventional methods (Vedantam et al., 2015). 306

To validate this evaluation method, we use the DOCCI dataset (Onoe et al., 2024), which contains human-annotated detailed image captions. Specifically, for each image in a subset of the dataset, we prepare the following four types of captions (details provided in Appendix B):

 Clean: The original caption (e.g., An indoor top-down view captures a white cat with black patches on a wooden floor, attempting to catch a large pale peacock feather flying above it.).

Object: An additional description of an object likely to exist in the image but not actually present is added to the Clean caption (*e.g.*, *An indoor top-down view captures a white cat with black patches on a wooden floor, attempting to catch a large pale peacock feather flying above it. A small red ball is rolling near the cat.*).

316 3. Attribution: Some object attributions in the Clean caption are modified to be inconsistent with the
 image (e.g., An indoor top-down view captures a white cat with black patches on a metal floor,
 attempting to catch a small dark peacock feather flying above it.).

4. Relation: Some relationships between objects in the Clean caption are altered to be inconsistent with the image (*e.g.*, *An indoor top-down view captures a white cat with black patches on a wooden floor, attempting to catch a large pale peacock feather flying below it.*).

We evaluate the four types of captions using various image caption evaluation methods (BLEU,
 ROUGE, METEOR, CIDEr, CLIP-S, RefCLIP-S, CLAIR, and ALOHa), including our own, to
 determine whether the hallucinations in the three modified types are reflected in the scores. For a fair

328

330

331

332

333

334

335 336

337

342

343

344

324



47. What is on the left side of the table in the photo?

- A) A lamp B) A blue decorative tree C) A vase D) A stack of books 48. What is the texture of the table in the foreground?
- 40. What is the texture of the table in the foreground r A) Smooth and shiny <u>B) Rough and rustic</u> (C) Soft and cushioned D) Metallic and cold 49. What is in the background of the photo, to the right side?
- A) A kitchen <u>B) A Christmas tree</u> C) A bookshelf D) A window with curtains
 50. What type of ornaments are on the triangular decoration on the table?
- A) Animal figurines <u>B) Christmas baubles</u> C) Miniature houses D) Candles

Figure 4: An example of our coverage evaluation data sample. The dataset consists of multiple-choice questions with four or fewer options. As demonstrated, the dataset includes questions with varying levels of granularity, ranging from broad to highly detailed. We have an LLM solve these problems using only the provided captions.

comparison, all methods requiring GPT (CLAIR, ALOHa, and ours) use GPT-40, and all methods requiring reference captions (BLEU, ROUGE, METEOR, CIDEr, CLAIR, ALOHa, and ours) use a separate set (Garg et al., 2024) of human-annotated captions (one reference caption per image).

Table 2 shows that existing metrics are unreliable for evaluating the factuality of detailed image 345 captions. Specifically, BLEU, ROUGE, METEOR, and CIDEr fail to account for hallucinations in 346 the scores and do not assign the highest score to the Clean captions. CLIP can only process up to 77 347 tokens and operates like a bag-of-words model (Yuksekgonul et al., 2023). This prevents CLIP-based 348 metrics from capturing the full content of detailed image captions, particularly missing Relation 349 hallucinations. ALOHa effectively addresses Object and Attribution hallucinations but fails to capture 350 Relation hallucinations due to its algorithmic limitations. CLAIR detects and reflects all three types 351 of hallucinations in the scores. However, CLAIR does not focus solely on factuality; instead, it allows 352 the GPT model to directly score each caption, applying the evaluation criteria implicitly defined by 353 the GPT model. In contrast, our evaluation method exclusively considers the factuality of the caption. 354 While it does not assign a perfect score to the Clean captions due to GPT-40's limitations in image understanding, it successfully assigns the highest score to Clean among the four caption sets. 355

356

Coverage. Even if an image caption contains only factual information, it would not be highly
 rated if it reflects only trivial aspects of the image. To assess the coverage of captioning models, we
 propose a QA-based metric and a benchmark dataset. Our coverage evaluation method is based on
 the assumption that *if an image caption fully captures the information in the image, visual questions about that image should be answerable by referencing the caption alone.*

Our goal is to evaluate detailed image captioning models. Therefore, the visual questions for evaluation must include a variety of detailed and nuanced questions about the images. Given the limitations of existing VQA datasets in this regard (Lu et al., 2022; Yin et al., 2023b; Li et al., 2023b; Yue et al., 2024), we construct a new VQA dataset. However, creating a new VQA dataset that includes a variety of detailed questions requires substantial labor. To reduce the associated costs, we follow the process outlined below to construct our dataset:

- 1. Generating more than 50 questions per image in the IIW-400 dataset using GPT-40.
- 2. Deduplicating the questions for each image using Sentence-BERT (Reimers & Gurevych, 2019).
- 370 3. Instructing human labelers to remove or revise questions that can be answered without specific 371 image information, or that are ambiguous or flawed, making them difficult to answer.
- 4. Annotating the correct answers to the remaining and revised questions by human labelers.

Our coverage evaluation dataset contains a total of 19,899 multiple-choice questions, with each image averaging 49.8 questions. Although we did not explicitly instruct the GPT model to generate detailed questions, it naturally includes them while generating a large number of questions. We present an example of our dataset in Figure 4. While our benchmark dataset can also be used to assess the visual understanding capabilities of MLLMs, we use it to evaluate the coverage of captioning models by having an LLM answer the questions based on the captions generated by those models.

381 382

397 398 399

400 401

402

Table 3: Effectiveness of our proposed method across various captioning models. In the V-FactER column, the LLM represents the decomposer and corrector, while the MLLM represents the fact-checker. Avg. denotes the average of the CLAIR, Factuality, and Coverage results.

Contionar	V	-FactER		Metri	ic	
Captioner	LLM	MLLM	CLAIR	Factuality	Coverage	Avg.
	-	-	68.8	59.9	47.9	58.9
LLaVA-NeXT-7B	LLaMA-3-8B	LLaVA-NeXT-7B	74.1	72.2	46.9	64.4
	GPT-4	LLaVA-NeXT-7B	74.6	73.4	46.2	64.7
	-	-	70.2	62.1	48.5	60.3
LLaVA-NeXT-13B	LLaMA-3-8B	LLaVA-NeXT-13B	75.5	77.9	45.8	66.4
	GPT-4	LLaVA-NeXT-13B	73.4	79.3	45.1	65.9
	-	-	74.9	65.5	48.2	62.9
InternVL-Chat-V1.5	LLaMA-3-8B	InternVL-Chat-V1.5	78.2	75.9	47.3	67.1
	GPT-4	InternVL-Chat-V1.5	77.8	75.7	47.3	66.9
	-	-	82.4	77.1	53.5	71.0
CDT 4V	LLaMA-3-8B	LLaVA-NeXT-7B	83.3	83.3	50.8	72.4
OF 1-4 V	LLaMA-3-8B	LLaVA-NeXT-13B	81.9	85.3	48.4	71.9
	LLaMA-3-8B	InternVL-Chat-V1.5	84.6	82.1	53.5	73.4

4 EXPERIMENTAL RESULTS AND DISCUSSION

4.1 EXPERIMENTAL SETUP

We adopt LLaVA-v1.5-7B, LLaVA-NeXT-7B, LLaVA-NeXT-13B, InternVL-Chat-V1.5, and GPT-4V
as the models for both captioning and V-FactER's fact-checking. We use LLaMA-3-8B (AI@Meta, 2024) or GPT-4 as the decomposer and corrector LLMs in V-FactER. Our experiments utilize the IIW-400 dataset, which contains 400 images, each accompanied by a highly detailed, hallucination-free caption. These high-quality reference captions enable precise evaluation of the captioning models.

We employ our proposed factuality and coverage evaluation methods, along with CLAIR, all of which use GPT-40 to evaluate the generated captions. To ensure robust evaluation and assess the recall potential of the captioning methods, we summarize the captions (Ge et al., 2024) generated from five different input prompts using LLaMA-3-8B. The only hyperparameter in V-FactER, π , is determined using a validation set composed of five images, their QAs, and reference captions. This validation set is constructed by sampling five examples from the DCI dataset (Urbanek et al., 2024). The prompt templates used in our experiments are provided in Appendix B.

415

4.2 IMPROVEMENT IN THE FACTUALITY OF CAPTIONING MODELS

⁴¹⁷Our proposed V-FactER exhibits a loose factuality-coverage trade-off depending on the hyperparameter π . Specifically, as π decreases, the threshold for determining factual propositions becomes stricter, leading to more propositions being identified for correction. Consequently, factuality increases while coverage decreases (an ablation study on π is provided in Appendix A). We first investigate whether V-FactER can enhance the factuality of various MLLMs while minimizing the reduction in coverage.

422 Table 3 demonstrates that V-FactER can significantly enhance the factuality of all tested MLLMs 423 while minimizing coverage loss. The substantial improvement in factuality, compared to the relatively 424 minor coverage loss in the captioning models, is also reflected in the increased CLAIR scores. Using 425 a more advanced LLM in V-FactER does not necessarily result in greater performance gains. 426 When applying V-FactER to the LLaVA and InternVL models, there is minimal difference between 427 the results obtained with LLaMA-3-8B and those with GPT-4. This suggests that the LLM's role in 428 V-FactER is relatively straightforward. V-FactER can improve detailed image captioning even 429 for the state-of-the-art MLLM, GPT-4V. It can significantly enhance factuality even when used with MLLMs far less capable than GPT-4V. However, in such cases, there is a considerable loss 430 in coverage, as many visual elements recognized by GPT-4V are identified as hallucinations by 431 V-FactER. With InternVL-Chat-V1.5, V-FactER maintains GPT-4V's coverage while improving

Table 4: Performance comparison between our proposed method and other methods regarding detailed image captioning. Base refers to the default image captioning of LLaVA-v1.5-7B without additional techniques.

Method	CLAIR	Factuality	Coverage	Avg.
Base	62.1	52.8	34.3	49.7
VCD (Leng et al., 2024)	59.7	44.6	39.3	47.9
OPERA (Huang et al., 2024)	59.1	53.0	34.1	48.7
Volcano (Lee et al., 2024)	63.9	53.7	37.7	51.7
LRV-Instruction (Liu et al., 2023a)	39.7	29.1	37.8	35.5
V-FactER (ours)	66.3	63.4	33.1	54.3

Table 5: Detailed image captioning and VOA performance of various MLLMs. OpenCompass (Duan et al., 2024) includes MMBench v1.1 (Liu et al., 2023b), MMStar (Chen et al., 2024a), MMMU val (Yue et al., 2024), MathVista (Lu et al., 2024), OCRBench (Liu et al., 2024d), AI2D (Kembhavi et al., 2016), HallusionBench (Guan et al., 2024), and MMVet (Yu et al., 2023). For POPE (Li et al., 2023b), we report the average F1 score across the three categories: adversarial, popular, and random. We report the sum of the perception and cognition scores for MME (Yin et al., 2023b). The best results for each metric are shown in **bold**.

Model	De	Detailed Image Captioning			Visual Question Answering			
1110401	CLAIR	Factuality	Coverage	Avg.	OpenCompass	MME	POPE	Avg.
InstructBLIP-7B	57.2	44.4	30.3	43.9	31.1	1391.4	86.1	38.4
LLaVA-v1.5-7B	61.1	56.3	30.5	49.3	36.9	1808.4	86.1	44.6
LLaVA-NeXT-7B	63.8	58.5	42.2	54.8	44.7	1769.1	87.5	50.8
LLaVA-NeXT-13E	64.5	62.8	43.0	56.8	47.6	1745.6	87.8	53.1
Idefics2-8B	58.1	85.2	13.4	52.2	53.0	1847.6	86.2	57.6
InternVL-Chat-V1.	5 72.4	67.6	46.0	62.0	61.7	2189.6	87.5	65.9
MiniCPM-V-2.6	73.1	68.9	43.6	61.9	65.2	2268.7	83.2	68.6
GPT-4V	82.4	78.6	52.6	71.2	63.5	2070.2	81.8	66.4

factuality. We additionally provide a qualitative comparison in Figure 5 between LLaVA-NeXT-7B with and without the application of V-FactER (referencing the first two rows of Table 3).

4.3 COMPARISON WITH OTHER METHODS

Various methods have been proposed to mitigate hallucinations in MLLMs, and they have primarily been validated on VQA and simple captioning benchmarks. We compare V-FactER with two recent decoding-based methods (VCD and OPERA), one corrector-based method (Volcano), and one training-based method (LRV-Instruction) from the perspective of detailed image captioning. All methods, except for LRV-Instruction, use LLaVA-v1.5-7B, while the LRV-Instruction method employs the MiniGPT-4 model (Zhu et al., 2023), as provided by its authors.

Table 4 shows that the VCD, OPERA, and LRV-Instruction methods are ineffective for detailed image captioning. Ironically, applying VCD significantly reduces the factuality of the LLaVA model while increasing coverage. Volcano yields only slight improvements in LLaVA's captions. However, V-FactER substantially enhances the factuality of the captioning model compared to the other methods. These results suggest that methods proposed to enhance MLLM factuality should be evaluated not only on tasks requiring short responses, such as VQA, but also on detailed image captioning tasks.

4.4 CONSISTENCY BETWEEN MLLM CAPTIONING AND VQA EVALUATION RESULTS

Currently, MLLM evaluations are primarily conducted on tasks that require only short responses, such as VQA tasks (Duan et al., 2024). However, to assess the potential of MLLMs in real-world applications, such as visual assistants, it is essential to evaluate their detailed image captioning abilities. The ranking of models used in our experiments, including LLaVA-v1.5-7B, LLaVA-NeXT-7B, LLaVA-NeXT-13B, InternVL-Chat-V1.5, and GPT-4V, is consistent across both our captioning evaluation results and widely used benchmarks like MMMU (Yue et al., 2024). However, for instance,



Figure 5: An example of a caption generated by V-FactER, with LLaVA-NeXT-7B as both the captioning and fact-checking model and LLaMA-3-8B as both the decomposer and corrector LLM.

some MLLMs may be optimized for VQA tasks that require only short responses, allowing them 502 to rank highly on common VQA benchmarks, yet their limited image captioning abilities could restrict their practical use. To investigate this, we evaluate the detailed image captioning capabilities of various MLLMs and examine whether their rankings are consistent with their rankings on widely used VQA benchmarks. We adopt InstructBLIP-7B (Dai et al., 2023), Idefics2-8B (Laurencon et al., 2024), and MiniCPM-V-2.6 (Yao et al., 2024) as additional MLLMs for the experiment.

Table 5 presents the evaluation results of MLLMs' responses to the prompt "Describe the given 508 image in a very detailed manner" as well as the performance of these models on various VQA tasks. 509 From these results, we observe that the performance of an MLLM on widely used benchmarks does 510 not necessarily reflect its capabilities in detailed image captioning. Specifically, Idefics2-8B ranks 511 mid-tier among the tested models in VQA tasks but falls into the lowest-performing group in terms 512 of detailed image captioning. Its high factuality but low coverage indicates that Idefics2-8B has been 513 trained to provide short and concise answers; this conclusion remains unchanged even when using 514 Idefics2-8B-Chatty (Laurençon et al., 2024). Despite being a relatively small model, MiniCPM-V-2.6 515 attracted attention by outperforming GPT-4V on benchmarks. However, our results show that the 516 model significantly underperforms GPT-4V in detailed image captioning. Additionally, we find that the factuality of the captions cannot be reliably predicted from the accuracy of MLLMs on POPE (Li 517 et al., 2023b), which was proposed to evaluate object hallucinations. 518

519 Based on these experimental results, we raise concerns about the current MLLM evaluations that 520 are centered around VQA tasks. We encourage the community to also evaluate MLLMs from the 521 perspective of detailed image captioning in order to showcase their full potential.

522 523 524

498

499 500 501

504

505

506

507

5 CONCLUSION

Detailed image captioning tasks are closely linked to critical applications, such as visual assistance for the impaired. Our research aims to assess and enhance the potential of MLLMs in these real-527 world contexts. We propose V-FactER, a method that improves detailed image captions through 528 the collaboration of a pre-trained MLLM and LLM. In addition, we introduce a framework and 529 benchmark dataset for evaluating the factuality and coverage of captioning models. Our experiments 530 validate the proposed evaluation framework and demonstrate that V-FactER significantly improves 531 the factuality of captioning models. We additionally present the following two key observations:

- 532 Methods designed to improve MLLM factuality, which have been validated primarily on VQA or 533 short captioning tasks, may be ineffective for detailed image captioning and can even reduce the 534 factuality of the backbone model's responses.
- High performance on commonly used VQA-centric benchmarks does not necessarily indicate that 536 the model will excel in detailed image captioning.
- These observations raise concerns about the current VQA-centric trend in MLLM evaluation. We 538 encourage the community to evaluate MLLMs and related algorithms not only on VQA tasks but also on detailed image captioning tasks to gain a more comprehensive understanding of their potential.

540 REPRODUCIBILITY STATEMENT

The prompt templates used in our proposed V-FactER are provided in Appendix B. The factuality and coverage evaluation codes are included in the supplementary material, along with a subset of our proposed benchmark dataset. The full dataset will be made publicly available soon.

References

542

543

544

545 546

547 548

549

550

551

552

553 554

555

556

558

559

560

565

566

567

576

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/ blob/main/MODEL_CARD.md.
- Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pp. 65–72, 2005.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- David Chan, Suzanne Petryk, Joseph Gonzalez, Trevor Darrell, and John Canny. Clair: Evaluating image captions with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 13638–13646, 2023.
 - Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*, 2024a.
- ⁵⁶⁸
 ⁵⁶⁹ Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024b.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong
 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning
 for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24185–24198, 2024c.
- Yiming Cui, Ziqing Yang, and Xin Yao. Efficient and effective text encoding for chinese llama and alpaca. *arXiv preprint arXiv:2304.08177*, 2023.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. InstructBLIP: Towards general-purpose vision-language models with instruction tuning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=vvoWPYqZJA.
- Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang
 Zang, Pan Zhang, Jiaqi Wang, Dahua Lin, and Kai Chen. Vlmevalkit: An open-source toolkit
 for evaluating large multi-modality models, 2024. URL https://arxiv.org/abs/2407.
 11691.
- Roopal Garg, Andrea Burns, Burcu Karagol Ayan, Yonatan Bitton, Ceslee Montgomery, Yasumasa
 Onoe, Andrew Bunner, Ranjay Krishna, Jason Baldridge, and Radu Soricut. Imageinwords:
 Unlocking hyper-detailed image descriptions. *arXiv preprint arXiv:2405.02793*, 2024.
- Yunhao Ge, Xiaohui Zeng, Jacob Samuel Huffman, Tsung-Yi Lin, Ming-Yu Liu, and Yin Cui. Visual fact checker: Enabling high-fidelity detailed caption generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14033–14042, 2024.

594	Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang
595	Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entan-
596	gled language hallucination and visual illusion in large vision-language models. In Proceedings of
597	the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 14375–14385, 2024.
598	
599	Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-
600	Methods in Natural Language Processing, pp. 7514, 7528, 2021
601	Memous în Matural Language Flocessing, pp. 7514–7528, 2021.
602	MD Zakir Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. A comprehensive
603	survey of deep learning for image captioning. ACM Computing Surveys (CsUR), 51(6):1–36, 2019.
604	
605	Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming
606	Zhang, and Nenghai Yu. Opera: Alleviating hallucination in multi-modal large language models
607	via over-trust penalty and retrospection-allocation. In <i>Proceedings of the IEEE/CVF Conference</i>
608	on Computer vision and Pattern Recognition, pp. 13418–13427, 2024.
609	Aniruddha Kembhayi, Mike Salvato, Eric Kolve, Minioon Seo, Hannaneh Hajishirzi, and Ali Farhadi
610	A diagram is worth a dozen images. In <i>Computer Vision–ECCV 2016: 14th European Confer-</i>
611	ence, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14, pp. 235–251.
612	Springer, 2016.
613	
614	Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building
615	vision-language models? arXiv preprint arXiv:2405.02246, 2024.
616	Seongyun Lee, Sue Park, Vongree Jo, and Minioon Seo, Volcano: Mitigating multimodal hallu
617	cination through self-feedback guided revision. In Proceedings of the 2024 Conference of the
618	North American Chapter of the Association for Computational Linguistics: Human Language
619	Technologies (Volume 1: Long Papers), pp. 391–404, 2024.
620	
621	Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing.
622	Mitigating object hallucinations in large vision-language models through visual contrastive decod-
623	ing. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition,
624	pp. 13872–13882, 2024.
625	Junnan Li Dongxu Li Silvio Savarese and Steven Hoi Blin-2: Bootstranning language-image
626	pre-training with frozen image encoders and large language models. In <i>International conference</i>
627	on machine learning, pp. 19730–19742. PMLR, 2023a.
628	0/11
629	Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object
630	hallucination in large vision-language models. arXiv preprint arXiv:2305.10355, 2023b.
631	Chin Vaw Lin Pouge: A package for automatic avaluation of summarized. In Text summarization
632	branches out pp. 74-81, 2004
633	<i>branches bui</i> , pp. 74-61, 2004.
634	Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Mitigating
635	hallucination in large multi-modal models via robust instruction tuning. In The Twelfth International
636	Conference on Learning Representations, 2023a.
637	
638	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.
630	Liava-next: improved reasoning, ocr, and world knowledge, January 2024a. URL https://
640	11ava = v1.g1(11ub, 10/b10g/2024 = 01 = 50 = 11ava = 11ext/.
641	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances in
642	neural information processing systems, 36, 2024b.
643	
64/	Shi Liu, Kecheng Zheng, and Wei Chen. Paying more attention to image: A training-free method for
645	alleviating hallucination in lvims. arXiv preprint arXiv:240/.217/1, 2024c.
646	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangho Zhao, Yike Yuan, Jiagi
647	Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? <i>arXiv preprint arXiv:2307.06281</i> , 2023b.

648 Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xucheng Yin, Cheng 649 lin Liu, Lianwen Jin, and Xiang Bai. On the hidden mystery of ocr in large multimodal models, 650 2024d. URL https://arxiv.org/abs/2305.07895. 651 Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, 652 Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for 653 science question answering. In The 36th Conference on Neural Information Processing Systems 654 (NeurIPS), 2022. 655 656 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, 657 Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In International Conference on Learning Representations 658 (ICLR), 2024. 659 Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and factuality 661 in abstractive summarization. In Proceedings of the 58th Annual Meeting of the Association for 662 Computational Linguistics, pp. 1906–1919, 2020. 663 Yasumasa Onoe, Sunayana Rane, Zachary Berger, Yonatan Bitton, Jaemin Cho, Roopal Garg, Alexan-664 der Ku, Zarana Parekh, Jordi Pont-Tuset, Garrett Tanzer, et al. Docci: Descriptions of connected 665 and contrasting images. arXiv preprint arXiv:2404.19753, 2024. 666 667 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic 668 evaluation of machine translation. In Proceedings of the 40th annual meeting of the Association 669 for Computational Linguistics, pp. 311–318, 2002. 670 Suzanne Petryk, David Chan, Anish Kachinthaya, Haodi Zou, John Canny, Joseph Gonzalez, and 671 Trevor Darrell. Aloha: A new measure for hallucination in captioning models. In Proceedings 672 of the 2024 Conference of the North American Chapter of the Association for Computational 673 Linguistics: Human Language Technologies (Volume 2: Short Papers), pp. 342–357, 2024. 674 675 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 676 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In International conference on machine learning, pp. 677 8748-8763. PMLR, 2021. 678 679 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. 680 In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing. 681 Association for Computational Linguistics, 11 2019. URL https://arxiv.org/abs/1908. 682 10084. 683 Jack Urbanek, Florian Bordes, Pietro Astolfi, Mary Williamson, Vasu Sharma, and Adriana Romero-684 Soriano. A picture is worth more than 77 text tokens: Evaluating clip-style models on dense 685 captions. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 686 pp. 26700-26709, 2024. 687 A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017. 688 689 Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image 690 description evaluation. In Proceedings of the IEEE conference on computer vision and pattern 691 recognition, pp. 4566-4575, 2015. 692 693 Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In Proceedings of the IEEE conference on computer vision and pattern 694 recognition, pp. 3156-3164, 2015. 696 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha 697 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In The Eleventh International Conference on Learning Representations, 2023. 699 Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich 700 Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual 701 attention. In International conference on machine learning, pp. 2048–2057. PMLR, 2015.

702 703 704	Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. <i>arXiv preprint arXiv:2408.01800</i> , 2024.
705 706 707	Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. <i>arXiv preprint arXiv:2306.13549</i> , 2023a.
708 709 710	Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. <i>arXiv preprint arXiv:2306.13549</i> , 2023b.
710 711 712 713	Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. Woodpecker: Hallucination correction for multimodal large language models. <i>arXiv preprint arXiv:2310.16045</i> , 2023c.
714 715 716	Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. <i>arXiv</i> preprint arXiv:2308.02490, 2023.
717 718 719 720 721 722	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In <i>Proceedings of CVPR</i> , 2024.
723 724 725	Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? In <i>The Eleventh International Conference on Learning Representations</i> , 2023.
726 727 728 729	Tianhang Zhang, Lin Qiu, Qipeng Guo, Cheng Deng, Yue Zhang, Zheng Zhang, Chenghu Zhou, Xinbing Wang, and Luoyi Fu. Enhancing uncertainty-based hallucination detection with stronger focus. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pp. 915–932, 2023.
730 731 732 733 734 735	Yukun Zhao, Lingyong Yan, Weiwei Sun, Guoliang Xing, Chong Meng, Shuaiqiang Wang, Zhicong Cheng, Zhaochun Ren, and Dawei Yin. Knowing what llms do not know: A simple yet effective self-detection method. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pp. 7044–7056, 2024.
736 737 738 739	Weihong Zhong, Xiaocheng Feng, Liang Zhao, Qiming Li, Lei Huang, Yuxuan Gu, Weitao Ma, Yuan Xu, and Bing Qin. Investigating and mitigating the multimodal hallucination snowballing in large vision-language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pp. 11991–12011, 2024.
740 741 742	Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and Huaxiu Yao. Analyzing and mitigating object hallucination in large vision-language models. In <i>The Twelfth International Conference on Learning Representations</i> , 2024.
743 744 745 746	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. <i>arXiv preprint arXiv:2304.10592</i> , 2023.
747 748 749	
750 751	
752	
754	

LLaMA-3-8B

LLaMA-3-8B

LLaMA-3-8B

A ABLATION STUDY

771

772

773

774

775

781 782

794

795

796

797

798

799

800 801

756

Captioner		V-FactER				Metric			
	LLM	MLLM	π	CLAIR	Factuality	Coverage			
	-	-	-	68.8	59.9	47.9			
LLaVA-NeXT-7B	LLaMA-3-8B	LLaVA-NeXT-7B	1.0	74.1	72.2	46.9			
	LLaMA-3-8B	LLaVA-NeXT-7B	0.5	73.6	76.9	43.7			
	LLaMA-3-8B	LLaVA-NeXT-7B	0.3	72.2	76.8	40.0			
	-	-	-	70.2	62.1	48.5			
LL WA NOVT 121	LLaMA-3-8B	LLaVA-NeXT-13B	1.0	75.5	77.9	45.8			
LLavA-NCAI-ISD	LLaMA-3-8B	LLaVA-NeXT-13B	0.5	74.8	79.9	42.1			
	LLaMA-3-8B	LLaVA-NeXT-13B	0.3	72.6	80.5	39.6			

74.9

78.2

79.0

77.7

1.0

0.5

0.3

65.5

75.9

78.8

81.7

48.2

47.3

46.0

42.5

Table 6: Effectiveness of our proposed method across various captioning models as a function of π . In the V-FactER column, the LLM represents the decomposer and corrector, while the MLLM represents the fact-checker.

Our proposed method features a single hyperparameter, π , which serves as the threshold for classifying atomic propositions as hallucinations or non-hallucinations. Table 6 presents the effects of V-FactER across various models as a function of π . The results reveal a loose trade-off between factuality and coverage depending on π . Specifically, in all tested settings, as π increases, factuality tends to decrease while coverage increases.

InternVL-Chat-V1.5

InternVL-Chat-V1.5

InternVL-Chat-V1.5

B PROMPT TEMPLATES

InternVL-Chat-V1.5

 prompt_1 = "Describe the given image in a very detailed manner."
<pre>prompt_2 = "Provide a detailed description of the specified image."</pre>
prompt_3 = "Elaborate on the details of the image provided."
prompt_4 = "Offer an in-depth description of the given image."
<pre>prompt_5 = "Thoroughly describe the features of the specified image."</pre>

Figure 6: The five prompt inputs used to generate captions in our experiments.

system:

I want to verify if the given CAPTION is accurate. To assist with this verification, decompose the given CAPTION into atomic propositions. All parts of the caption must be broken down into propositions. The outputs should follow the following format:'1. proposition one\n2. proposition two\n3. proposition three'. For example, break down 'He is tall, thin, and pale' into '1. He is tall.\n2. He is thin.\n3. He is pale.'

user: CAPTION: {caption}

Figure 7: The prompt input for LLaMA-3-8B serving as the decomposer.

- 805
- 806
- 807
- 809

⁸⁰⁴

813

814

815

816

817 818

819 820 821

822

823

824

825

827

828 829

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848 849

850 851 852 system:

user:

user:

I want to create a caption that includes only facts. Please help me correct the given caption. The given caption contain things that are not true. Based on the given FACTS and NON-FACTS remove the non-factual elements from the caption. Place the revised caption between '###'.

Caption: {caption}\nFACTS:\n{Non-hallucinations among the atomic propositions}\nNON-FACTS:\n{Hallucinations among the atomic propositions}

Figure 8: The prompt input for LLaMA-3-8B serving as the corrector.

system: This is a hard problem. Carefully summarize in ONE detailed caption based on the following 5 captions by different people describing the same image. Be sure to describe everything, and avoid hallucination. Provide the detailed caption in the format '### {Detailed caption} ###'.

Caption 1: {caption 1st}\n Caption 2: {caption 2nd}\n Caption 3: {caption 3rd}\n Caption 4: {caption 4th}\n Caption 5: {caption 5th}\n

Figure 9: The prompt input for LLaMA-3-8B serving as the summerizer. We use the prompt employed in the work of Ge et al. (2024).

if hallucination == "Object"

prompt_sys = "I want to inject incorrect information into the caption of the given photo. Your role is to modify about THREE words from the latter part of the given caption that describe the attributes of the objects so that they do not match the photo." elif hallucination == "Attribution" prompt_sys = "I want to inject incorrect information into the caption of the given photo. Your role is to imagine an object that isn't actually in the image but could plausibly be there, and add a very brief part about it to the caption so that they do not match the photo." elif hallucination == "Relation":

prompt_sys = "I want to inject incorrect information into the caption of the given photo. Your role is to change the spatial relationships between the objects so that they do not match the photo. For example, change 'A person is standing to the right of the car' to 'A person is standing to the left of the car.' Do not change anything other than the spatial relationships between the objects."

system:
{prompt_sys}

svstem:

user

user (presented with the image):
Caption: {caption}

together with \";\" as separation."

Figure 10: The prompt input for GPT-40 used to create the meta-evaluation dataset of Table 2.

"I want to use an object detector to check the correctness of an image caption obtained by an image caption model. Can you help to parse the given CAPTION and list all objects that could

be detected with an object detection model in the image? Please only list the object name and

ignore the description. Please use the name in the CAPTION as it is. Please concatenate them

CAPTION: {caption}

Figure 11: The prompt input for GPT-4 used to create the dataset of Figure 1. We use the prompt employed in the work of Ge et al. (2024).

862 863

859