000 MITIGATING MULTIMODAL HALLUCINATIONS VIA 001 GRADIENT-BASED SELF-REFLECTION 002 003

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

032

033

034

037

040

041

042

043

044

045

Paper under double-blind review

ABSTRACT

Hallucination in Multimodal Large Language Models (MLLMs) occurs when inaccurate text-visual alignments are generated, posing a major challenge for reliable model output. Previous studies have identified three primary biases as major causes of hallucinations: text-visual bias (over-reliance on text over visual details), co-occurrence bias (misleading object correlations), and long-term bias (increased hallucinations in later stages of long sequences). Existing hallucination mitigation methods often rely on visual grounding, which requires additional resources such as scoring systems using another MLLM, and still fail to fully address all biases, particularly co-occurrence bias in visual inputs. We propose Gradient-based Influence-Aware Contrastive Decoding (GACD) to explicitly and jointly balance these biases, thereby mitigating hallucinations. To quantify these biases at the individual sample level, we introduce 'token influence'. Since biases are rooted in the training data and become embedded in pre-trained MLLMs, we derive token influence through self-reflection by calculating the gradients from output predictions to input tokens. Notably, GACD is the first approach capable of fully addressing co-occurrence bias without relying on extra resources or any form of tuning. Extensive experiments demonstrate GACD's effectiveness in reducing hallucinations and improving MLLM performance, achieving new stateof-the-art results while providing insights into the visual perception capabilities of these models.

Cause of hallucination text-visual bias co-occurrence bias long-term bias

the top right corner.

LLaVA v1.5: The image features a dining table with a LLaVA v1.5-GACD: The image features a dining plate of food, a fork, and a knife placed on it. The plate table with a plate with remnants of food and a fork on appears to be half-eaten, and there is a glass of beer next it. The plate has a few forkfuls of food, and there is a to it. The table is set with various utensils, including a glass nearby. The table also has a bottle, possibly spoon and another fork, as well as a knife. In addition to containing a beverage, and a cup. In addition to the the main dining setup, there are two bowls on the table, main items, there are several utensils on the table. one located near the top left corner and the other towards including a spoon, and another fork < 7% Early Stop



046 Figure 1: An illustration of our proposed GACD method and its effectiveness in addressing hal-047 lucinations caused by three key biases: text-visual bias, co-occurrence bias, and long-term bias. 048 GACD measures token influence via gradients to uncover biases embedded in pre-trained MLLMs. Specifically, GACD-VA balances text-visual bias by increasing the visual influence to match the influence from the prompt or the previous output. GACD-CR detects distracting-object-related visual 051 tokens and amplifies the influence of visual tokens unrelated to these distracting objects, effectively mitigating co-occurrence hallucinations. Additionally, an early stopping strategy based on visual 052 influence is introduced to prevent long-term hallucinations.

⁵⁴ 1 INTRODUCTION

055 056

Multimodal Large Language Models (MLLMs) have recently gained significant momentum due to their ability to process and generate text from various data modalities, including images and videos. Despite their impressive capabilities, these models face a significant challenge: hallucination, where the generated text content is not grounded in the visual inputs.

Hallucinations in MLLMs arise from three key biases extensively discussed in the literature Li et al. 061 (2023); Kang & Choi (2023); Kim et al. (2024): text-visual bias, co-occurrence bias, and long-term 062 bias. Text-visual bias arises when a model relies too heavily on textual input, neglecting visual 063 information. Co-occurrence bias results from statistical correlations in training data, leading mod-064 els to predict objects based on frequent co-occurrence, even when absent from visual inputs. For 065 instance, in a LLaVA-1.5 generated description (see Fig. 1), 'knife' and 'beer' are described, while 066 'fork' and 'glass' appear in the input image. Long-term bias, highlighted in yellow in Fig. 1 and 067 analyzed as 'position factors' in Zhou et al. (2023); Favero et al. (2024), refers to an increased like-068 lihood of hallucinations in later stages of long-form sequences, often caused by earlier predictions 069 overly influencing subsequent content generation.

Existing methods mostly focus on mitigating visual-text bias by grounding predictions in visual in-071 puts. These methods typically require additional resources to check predictions, such as scoring 072 systems that involve another MLLM prone to hallucination Xing et al. (2024); Deng et al. (2024); 073 Radford et al. (2021), or use contrastive decoding Leng et al. (2023); Favero et al. (2024) without 074 precise mechanisms to balance bias. On the other hand, because these approaches ground predic-075 tions on visual inputs alone, they fail to eliminate distractions caused by irrelevant or misleading 076 visual information, thus co-occurrence bias remains a problem. Long-term bias is usually addressed in the literature by incorporating sequence length as a model parameter Zhou et al. (2023); Favero 077 et al. (2024), capturing only a coarse likelihood of hallucination based solely on sequence length. Among all these existing methods, bias mitigation often depends on fixed hyperparameters uni-079 formly applied across all samples, requiring costly tuning and would underperform considering the 080 fact that biases behave differently for each single sample. There are also methods that mitigate bias 081 by training Chen et al. (2023); Jiang et al. (2024); Yue et al. (2024); Ben-Kish et al. (2023); Sun et al. (2023) and reinstructing the model through fine-tuning, which requires a huge cost on data 083 collection and model training. Therefore, how to jointly and reliably analyze and address all these 084 biases free from additional training or inputs remains an unresolved and challenging issue. 085

In this work, we propose Gradient-based Influence-Aware Contrastive Decoding (GACD) to solve 086 the problem. We first introduce a mathematical examination termed 'token influence' to reflect the 087 sensitivity between output and input tokens (visual, prompt, previous output in Fig. 1), using gradi-088 ent norms. This sample-wise, gradient-based influence reveals the local linear behavior of MLLMs, 089 showing that small changes in input lead to proportional changes in output, thereby exposing inherent biases embedded in the model. We then mitigate co-occurrence hallucinations, the first one 091 addressing this type of hallucination to the authors' knowledge, by amplifying the influence of visual 092 tokens unrelated to distractions (GACD-CR), and reduce visual-text hallucinations by amplifying all 093 visual input (GACD-VA). Additionally, to address long-term bias, we introduce a precise stopping criterion that halts generation when the visual influence ratio falls below a set threshold. 094

095 To be specific, we find that hallucinations are largely driven by the ratio of visual token influence 096 relative to the influence of all input tokens. This ratio serves as the underlying factor linking text-097 visual bias, co-occurrence bias, and long-term bias, thereby explaining how these biases contribute 098 to hallucinations. To address co-occurrence bias, GACD-CR identifies the visual tokens associated with the distracting object and labels it as a distracting-object-related token. For example, "fork" is a distracting object, and its related visual tokens are detected in Fig. 1. We categorize visual to-100 kens into two groups: those related to the distracting object and those unrelated. We then employ 101 influence-aware contrastive decoding to explicitly enhance the influence of essential visual tokens 102 and reduce the impact from distracting visuals to outputs. In GACD-VA, we amplify the influence 103 of all visual tokens (only those unrelated to the distracting object when combined with CR). This is 104 done to ensure that the influence of these tokens matches the maximum influence of other compo-105 nents, establishing it as the dominant influence to mitigate the text-visual bias, as illustrated in the 106 lower-right part of Fig. 1. 107

We summarize our contributions as follows:

tracting objects present within the visual inputs.

LLaVA-QA90 Liu et al. (2024b) across multiple MLLMs.

109

108

- 111 112
- 113 114
- 115 116
- 117
- 117 118
- 119

121 122

120

2 RELATED WORK

Hallucination and Bias. Wang & Sennrich (2020) demonstrated that hallucinations are more pro-123 nounced in out-of-domain test sets. As noted by Tonmoy et al. (2024); Li et al. (2023); Fu et al. 124 (2024), hallucinations are closely related to biases, particularly text-visual and co-occurrence bias. 125 Long-term bias, inherited from LLMs, has been studied in Favero et al. (2024); Zhou et al. (2023). 126 Existing methods Li et al. (2023); Fu et al. (2024); Kim et al. (2024) typically report only overall 127 statistics, lacking a mathematical sample-wise bias analysis. This distinction is crucial, as biases 128 can vary from case to case. We argue that biases are embedded in the parameters of pre-trained 129 MLLMs, as they inherently reflect the biases present in the training dataset. Therefore, analyzing 130 sample-wise bias through self-reflection offers a straightforward approach.

• We propose GACD, a novel hallucination mitigation method that enhances the influence of essential visual tokens and jointly addresses text-visual, co-occurrence, and long-term bias

via gradient-based self-reflection, without requiring external resources, training, or tuning.

• GACD is the first method capable of fully addressing co-occurrence bias, even when dis-

• We propose token influence and formulate bias problems of MLLMs through mathematical

GACD significantly improves baseline MLLMs, establishing new SOTA results. It achieves

up to a 6.2 F1 boost on POPE Li et al. (2023) and up to a 2.97 accuracy increase on the

expressions, enabling a detailed understanding of biases at the individual sample level.

131 Hallucination Mitigation. Existing hallucination mitigation methods can be grouped into two cate-132 gories: training-related and inference-related methods. Training-related methods Chen et al. (2023); 133 Jiang et al. (2024); Yue et al. (2024); Ben-Kish et al. (2023); Sun et al. (2023) are expensive, re-134 quiring access to training data and specialized modifications to that data to effectively mitigate 135 hallucinations. Inference-related methods are further divided into revision methods and guided or 136 contrastive decoding methods. Revision methods Xing et al. (2024); Deng et al. (2024); Radford 137 et al. (2021); Zhai et al. (2024) often use additional scoring systems, typically involving another 138 MLLM, which itself may hallucinate. In contrast, guided or contrastive decoding methods Leng et al. (2023); Zhao et al. (2024); Favero et al. (2024), adjust logits during inference without extra 139 training, making them easily integrable as add-on modules to any MLLM. Our approach falls into 140 this category. However, existing decoding methods still rely on extra resources. For instance, Leng 141 et al. (2023) utilizes noise images for contrastive decoding, while Zhao et al. (2024) employs a visual 142 decoder for guided decoding. Furthermore, their decoding weights are fixed using hyperparameters 143 Leng et al. (2023); Zhao et al. (2024) or adjusted based on the length of the generated text Favero 144 et al. (2024). In contrast, our method does not rely on any external resources and tuning. Instead, 145 it employs gradient-based influence analysis to calculate precise, sample-specific decoding weights. 146 This approach accurately balances multiple biases, leading to improved hallucination mitigation.

147 148

149 150

3 BIASES IN MULTIMODAL LARGE LANGUAGE MODELS

In this section, we analyze three key biases that underlie hallucinations in MLLMs: text-visual bias,
 co-occurrence bias, and long-term bias. We begin by introducing MLLMs, followed by a discussion
 on measuring the influence of input tokens, and conclude with an in-depth analysis of each bias.

Introduction to MLLMs. Consider a sequence $\mathbf{t}^p = [t_1^p, \dots, t_N^p]$ as the input prompt, where t_n^p ($1 \le n \le N$) is a prompt token from a predefined vocabulary. Visual tokens $\mathbf{t}^v = [t_1^v, \dots, t_S^v]$ are extracted from visual inputs V, using $\mathbf{t}^v = \mathcal{E}_{\tau}(V)$, where $\mathcal{E}_{\tau}(\cdot)$ is a visual encoder. These tokens are mapped to the token space via linear projection. The MLLMs generate the response sequence $\mathbf{y} = [y_1, \dots, y_M]$ using the logit generation function $\mathcal{F}_{\theta}(\cdot)$ and the softmax function $\sigma(\cdot)^{-1}$ as:

159
160
$$\mathbf{y} = \sigma(\mathcal{F}_{\theta}(\mathbf{t}^v, \mathbf{t}^p)),$$
 (1)

¹⁶¹

¹The output probability of the softmax function is interpreted as confidence in this paper.

where y_m denotes an individual output token for $1 \le m \le M$. The conditional probability distribution $p(\mathbf{y}|\mathcal{F}_{\theta}, \mathbf{t}^p, \mathbf{t}^v)$ can therefore be expressed as:

$$p(\mathbf{y}|\mathcal{F}_{\theta}, \mathbf{t}^{v}, \mathbf{t}^{p}) = \prod_{m=1}^{M} p(y_{m}|\mathcal{F}_{\theta}, \mathbf{t}^{v}, \mathbf{t}^{p}, \mathbf{y}_{< m}),$$
(2)

where $\mathbf{y}_{< m} = [y_1, \dots, y_{m-1}]$ for m > 1 and is empty for m = 1. The pre-trained MLLMs (θ^* and τ^*) are trained to minimize the overall loss of the entire training data:

$$\tau^{\star}, \theta^{\star} = \arg\min_{\tau, \theta} \mathcal{L}(\sigma(\mathcal{F}_{\theta}(\mathcal{E}_{\tau}(V), \mathbf{t}^{p})), \mathbf{y}^{gt}),$$
(3)

where $\mathcal{L}(\cdot)$ is the loss function of MLLMs, and \mathbf{y}^{gt} denotes the ground truth outputs. During training, MLLMs learn not only factual knowledge but also statistical correlations and biases present in the data, which become embedded in the model's parameters. As a result, these biases are captured within the pre-trained parameters. We aim to analyze these biases by analyzing the parameters τ^* .

Token Influence Measurement. To quantify token influence, we analyze the behavior of the model function $\mathcal{F}_{\theta^*}(\mathbf{t}^v, \mathbf{t}^p)$ for prediction at the m^{th} output, focusing on the local region around the visual embedding $\mathbf{t}^{v(0)}$ and the text prompt $\mathbf{t}^{p(0)}$. The linear behavior in the local area is likely to contain information about the preference. A straightforward approach is to consider the first order Taylor expansion ² of \mathcal{F}_{θ^*} w.r.t $\mathbf{t}^v, \mathbf{t}^p$ and, all the previous output tokens \mathbf{y}_i with i < m:

$$\mathcal{F}_{\theta^*}(\mathbf{t}^v, \mathbf{t}^p)_m \approx \sum_{s=1}^{S} \mathbf{g}_{ms}^v \cdot t_s^v + \sum_{n=1}^{N} \mathbf{g}_{mn}^p \cdot t_n^p + \sum_{i=1}^{m-1} \mathbf{g}_{mi}^y \cdot y_i + Const,$$
(4)

where:

165 166 167

170 171

182

183

185

187

196 197

200

205 206

$$\mathbf{g}_{ms}^{v} := \left. \frac{\partial (\mathcal{F}_{\theta^*})_m}{\partial t_s^{v}} \right|_{\mathbf{t}^{v} = \mathbf{t}^{v(0)}}, \quad \mathbf{g}_{mn}^{p} := \left. \frac{\partial (\mathcal{F}_{\theta^*})_m}{\partial t_n^{p}} \right|_{\mathbf{t}^{p} = \mathbf{t}^{p(0)}}, \quad \mathbf{g}_{mi}^{y} := \frac{\partial (\mathcal{F}_{\theta^*})_m}{\partial y_i} \tag{5}$$

are the first order gradients terms on visual tokens \mathbf{t}^v , text prompt \mathbf{t}^p and previous outputs y_i , Const denotes all other terms that are constant w.r.t the \mathbf{t}^v , \mathbf{t}^p .

In this sense, $\mathcal{F}_{\theta^*}(\mathbf{t}^v, \mathbf{t}^p)$ can be seen as a linear classifier w.r.t the visual tokens, text prompt and the previous outputs around a certain data point $(\mathbf{t}^{v(0)}, \mathbf{t}^{p(0)})$. The influence of tokens is indicated by the weight norm of the linear classifier, specifically, the gradient norm. We choose the Manhattan (L1) norm for our analysis and experiments because its sparsity clearly identifies the contribution of individual tokens. Consequently, the influence of each token can be represented as $|\mathbf{g}_{m}^{v/p/y}|$, and the overall influence of visual tokens, text prompt, and previous outputs can be represented as:

$$\mathbb{I}^{v} := \sum_{s=1}^{S} |\mathbf{g}_{ms}^{v}|, \quad \mathbb{I}^{p} := \sum_{n=1}^{N} |\mathbf{g}_{mn}^{p}|, \quad \mathbb{I}^{y} := \sum_{i=1}^{m-1} |\mathbf{g}_{mi}^{y}|.$$
(6)

This token influence expression allows us to conduct a mathematical analysis of sample-wise biases.

Text-visual Bias occurs when the influence of the text prompt tokens significantly outweighs that of the visual tokens, which can be expressed as $\mathbb{I}^p \gg \mathbb{I}^v$ or $\mathbb{I}^y \gg \mathbb{I}^v$, according to our proposed components influence. We begin our analysis of this bias with an object existence Visual Question Answering (VQA) task, involving 'yes/no' predictions, by calculating the overall influence ratio:

$$r^{v} = \frac{\mathbf{I}^{v}}{\mathbf{I}^{v} + \mathbf{I}^{p} + \mathbf{I}^{y}}, \quad r^{p} = \frac{\mathbf{I}^{p}}{\mathbf{I}^{v} + \mathbf{I}^{p} + \mathbf{I}^{y}}, \quad r^{y} = \frac{\mathbf{I}^{y}}{\mathbf{I}^{v} + \mathbf{I}^{p} + \mathbf{I}^{y}},$$
 (7)

on the POPE dataset Li et al. (2023). In this case, I^y is always zero, as the answer is a single word based on the question. The blue bars in Fig. 5a, which display visual influence ratio across LLaVA-v1.5, InstructBLIP, and mPLUG-Owl2, show that the prompt input exerts more influence than the visual tokens. This phenomenon is common in MLLMs and can be attributed to their training process, where multimodal features are aligned with language tokens after extensive textbased pre-training, causing language components to dominate the decision-making process.

We further observe text-visual bias in the open-ended image captioning task (Fig. 2) and the VQA task (Fig. 1 in the Appendix). Hallucinated predictions, highlighted in red, consistently exhibit a low

²For further details on the first-order Taylor expansion, please refer to the Appendix.







Figure 2: Comparison of component influence ratios in image captioning with and without GACD-VA. (Left) Captions generated using LLaVA-v1.5 show that as sequence length increases, the influence of previous output tokens continues to rise, diverging from the influence of visual tokens and prompt tokens. Hallucinated predictions, highlighted in red, are characterized by a low visual influence. (Right) GACD-VA amplifies the influence of visual tokens, thereby increasing confidence and reducing hallucinations. Punctuation marks and suffixes exhibit a low visual influence ratio.

240 visual influence ratio. As more tokens are generated, the influence of both visual and prompt tokens 241 diminishes, while the contribution of previous output tokens increases. This behavior aligns with 242 findings from previous studies Zhou et al. (2023); Favero et al. (2024) and reflects the tendency of 243 MLLMs to prioritize more recent elements in the sequence over earlier inputs. Additionally, the nature of the output token affects the degree of visual influence. For instance, punctuation marks 244 (such as '.') or suffixes (such as 'ably' in 'comfortably') tend to have a lower visual influence 245 ratio. This is intuitive, as these tokens rely more on linguistic context and are less dependent on 246 visual information. This observation underscores the value of hallucination mitigation methods 247 based on influence analysis, with GACD-VA (Sec. 4.1) offering more meaningful insights compared 248 to approaches that focus solely on the length of the generated sequence. 249

Co-occurrence Bias occurs when models overly rely on associations between elements that frequently appear together in the training data. Fig. 3a illustrates an example of co-occurrence hallucination, where mPLUG-Owl2 with GACD-VA (mPLUGOwl-VA) incorrectly predicts 'dining table' in the presence of 'chair' in the image, demonstrating that merely amplifying visual influence is insufficient to mitigate such hallucinations. We analyze the influence of visual tokens on the predic-



(a) Co-occurrence francemation

Figure 3: (a) Co-occurrence hallucination of 'dining table' in the presence of a 'chair'. (b) Visualization of individual visual token influence for the left example shows that the token with the highest influence on 'chair' also has the highest influence on 'table'. (c) Summary statistics for 100 chaironly and 100 table-only images, including hallucination rate and the percentage of cases where both objects share the same most influential visual token. GACD-CR successfully reduces both metrics.

268 269

216

217

218

219

220

221 222

228

229 230 231

232

233

tions of the distracting object ('chair'), the hallucinated object ('table'), and unrelated correct object ('bag'). Fig. 3b shows individual visual token influences: $|\mathbf{g}_{m_cs}^v|_{y_{m_c}=chair}$, $|\mathbf{g}_{m_ts}^v|_{y_{m_t}=table}$, and

 $|\mathbf{g}_{m_bs}^v|_{y_{m_b}=bag}$, revealing that the visual token most influential for the hallucinated 'table' also sig-270 271 nificantly contributes to 'chair', whereas the unrelated 'bag' is influenced by different visual tokens. 272

To further investigate, we collected 100 chair-only and 100 table-only images from the MSCOCO 273 Lin et al. (2014) evaluation dataset. Experimental results, shown in the blue bar in Fig. 3c, reveal 274 that in the presence of either a 'chair' or 'table' in the image, the hallucination rate for the other 275 object is 23.5%, with a shared most influential visual token rate of 31.9% in these hallucinations, 276 confirming our observation. This finding also inspired the development of GACD-CR (Sec. 4.2). 277

278 Long-term bias typically arises when generating long outputs, as hallucinations become more likely the further the output drifts from the visual context. This trend is evident in the LLaVA-v1.5 results 279 shown in Fig. 2, where hallucinations like 'bottle', 'cup', and 'spoon' occur as the influence of visual 280 tokens r^{v} decreases with increasing m. Based on this, we also propose a stopping criterion. 281

282 283

284

298

299

300

301

302

303

304

305

306

307

308

309

310

311

313 314

315

322

323

4 SELF-REFLECTIVE CONTRASTIVE DECODING

285 Building on the bias analysis, we propose a novel approach, GACD, as outlined in Fig. 4, which 286 uses self-reflection to achieve token influence awareness and adjust logits during decoding. This 287 approach follows the principle of contrastive decoding, aiming to enlarge the Kullback-Leibler (KL) divergence ³ between the contrastive logits distribution, generated by a subset of input tokens, and 288 the joint logits distribution, generated by all input tokens. This increase in KL divergence enhances 289 the influence of tokens outside the subset. The contrastive decoding weight, α , is calculated based 290 on the influence of visual, prompt, and previous outputs on both the original and contrastive logits, 291 ensuring that the adjusted logits receive a balanced visual influence. This calculation relies on 292 a local linear assumption of MLLM functions. While not strictly precise, it offers significantly 293 greater accuracy compared to using a fixed hyperparameter or relying solely on sequence length. 294 Additionally, we propose a stopping criterion based on the visual influence ratio r^{v} . If the visual 295 ratio for the sentence-starting token falls below a defined threshold, early stopping is triggered to 296 halt further output generation. 297



Figure 4: Overview of the proposed GACD method. It uses contrastive decoding to amplify the influence of essential visual tokens by generating contrastive logits from all input tokens except the essential visual ones. In GACD-VA (Sec. 4.1), contrastive logits are generated from prompt and previous output tokens to amplify all visual tokens. In GACD-CR (Sec. 4.2), visual tokens related to distracting objects are first detected, and then, along with prompt and previous output tokens, are 312 used to generate contrastive logits, thereby amplifying the influence of non-distracting visual tokens.

4.1 VISUAL-TOKEN AMPLIFICATION

316 We aim to increase the influence of visual tokens, ensuring alignment with both the prompt and the 317 previous output tokens. In VQA tasks, it is crucial for visual token influence to match the question 318 prompt to ensure visually grounded responses. Likewise, in open-ended generation, maintaining visual token influence in balance with previous outputs is essential to prevent visual forgetting. To 319 achieve this, our method adjusts the logits by incorporating contrastive logits generated from the 320 prompt and previous output tokens: 321

$$\hat{\mathcal{F}}_{\theta^*}(\mathbf{t}^v, \mathbf{t}^p)_m = (1 + \alpha_m) \mathcal{F}_{\theta^*}(\mathbf{t}^v, \mathbf{t}^p)_m - \alpha_m \mathcal{F}_{\theta^*}(\mathbf{t}^p)_m, \tag{8}$$

³A detailed explanation of why contrastive decoding enlarges KL divergence is provided in the Appendix.

where α_m is calculated by considering the influence of each component and explicitly matching the dominant influence from either the prompt or previous outputs:

$$\alpha_m = \frac{\mathbf{I}^{\star} - \mathbf{I}^v}{\mathbf{I}^v + \hat{\mathbf{I}}^{\star} - \mathbf{I}^{\star}}, \quad \text{where } \star = \begin{cases} p & \text{if } \mathbf{I}^p \ge \mathbf{I}^y \\ y & \text{otherwise} \end{cases}, \tag{9}$$

where $\hat{1}^p := \sum_{n=1}^N \left| \frac{\partial (\mathcal{F}_{\theta^*}(\mathbf{t}^p))_m}{\partial t_n^p} \right|$ and $\hat{1}^y := \sum_{i=1}^{m-1} \left| \frac{\partial (\mathcal{F}_{\theta^*}(\mathbf{t}^p))_m}{\partial y_i} \right|$ represent the influence of the prompt and previous output, respectively, on the contrastive logits.

Furthermore, unlike existing contrastive decoding methods Zhou et al. (2023); Favero et al. (2024); Leng et al. (2023), which rely on adaptive plausibility constraints such as the prediction confidence and require experiments to determine the optimal threshold, our approach operates as long as the prediction confidence is below 100%. Instead, our method explicitly imposes a constraint to prevent the influence of the prompt from becoming negative:

$$\alpha_m \le \frac{\mathbf{I}^p}{\hat{\mathbf{I}}^p - \mathbf{I}^p}, \quad \text{if } \hat{\mathbf{I}}^p > \mathbf{I}^p.$$
(10)

4.2 **CO-OCCURRENCE REDUCTION**

To reduce co-occurrence hallucination, we first identify visual tokens related to distracting objects by analyzing those with the highest influence, then mitigate their impact in subsequent predictions. A mask is applied to flag these distracting-object-related visual tokens, focusing on object (noun) output tokens. When a noun is predicted, the corresponding visual token with the highest influence is flagged. The mask \mathcal{M}_{ms}^v aggregates visual tokens related to all previously mentioned objects:

$$\mathcal{M}_{ms}^{v} = \mathcal{M}_{1s}^{v} \lor \mathcal{M}_{2s}^{v} \lor \cdots \lor \mathcal{M}_{(m-1)s}^{v}, \quad \mathcal{M}_{is}^{v} = \begin{cases} 1 & \text{if } \left| \mathbf{g}_{is}^{vc} \right| \ge \left| \mathbf{g}_{i(1...S)}^{vc} \right|, y_{i} \text{ is noun} \\ 0 & \text{otherwise} \end{cases}, \quad (11)$$

> where c in g_{is}^{vc} indicates that the gradient is specifically from the predicted object class, rather than from the overall predicted logits across the entire vocabulary classes.

> During prediction, if the output is a noun, contrastive logits are generated from both text and distracting-object-related visual tokens, modifying equation 8 as follows:

$$\tilde{\mathcal{F}}_{\theta^*}(\mathbf{t}^v, \mathbf{t}^p)_m = (1 + \tilde{\alpha}_m) \mathcal{F}_{\theta^*}(\mathbf{t}^v, \mathbf{t}^p)_m - \tilde{\alpha}_m \mathcal{F}_{\theta^*}(\mathbf{t}^v \times \mathcal{M}_m^v, \mathbf{t}^p)_m,$$
(12)

where $\times \mathcal{M}^v$ indicates the mask selection, and $\tilde{\alpha}_m$ controls the KL divergence between logits con-ditioned on only distracting-object-related visual tokens and those conditioned on all visual tokens. We modify equation 9 as follows:

$$\tilde{\alpha}_m = \frac{\mathtt{I}^{\star} - \mathtt{I}^v}{\mathtt{I}^v - \tilde{\mathtt{I}}^v + \tilde{\mathtt{I}}^{\star} - \mathtt{I}^{\star}}, \quad \star = \begin{cases} p & \text{if } \mathtt{I}^p \ge \mathtt{I}^y \\ y & \text{otherwise} \end{cases},$$
(13)

where $\tilde{\mathbf{I}}^v := \sum_{s=1}^{S} \left| \frac{\partial (\mathcal{F}_{\theta^*}(\mathbf{t}^v \times \mathcal{M}_{ms}^v, \mathbf{t}^p))_m}{\partial t_s^v} \right|$, $\tilde{\mathbf{I}}^p := \sum_{n=1}^{N} \left| \frac{\partial (\mathcal{F}_{\theta^*}(\mathbf{t}^v \times \mathcal{M}_{ms}^v, \mathbf{t}^p))_m}{\partial t_n^p} \right|$ and $\tilde{\mathbf{I}}^y := \sum_{i=1}^{m-1} \left| \frac{\partial (\mathcal{F}_{\theta^*}(\mathbf{t}^v \times \mathcal{M}_{ms}^v, \mathbf{t}^p))_m}{\partial y_i} \right|$ represent the influence of visual, prompt, and previous output tokens, respectively, on the contractive locity

respectively, on the contrastive logits. Since distracting-object-related visual tokens are included to generate the contrastive logits, along-

side prompt tokens, we ensure that the influence of these visual tokens does not become negative like that of the prompt tokens. Accordingly, we modify the constraint in equation 10 as follows:

$$\tilde{\alpha}_m \le \min\left(\frac{\mathbf{I}^d}{\tilde{\mathbf{I}}^v - \mathbf{I}^d}, \frac{\mathbf{I}^p}{\tilde{\mathbf{I}}^p - \mathbf{I}^p}\right), \quad \text{if } \tilde{\mathbf{I}}^v > \mathbf{I}^d, \quad \tilde{\mathbf{I}}^p > \mathbf{I}^p, \tag{14}$$

where $\mathbf{I}^d = \sum_{s=1}^{S} |\mathbf{g}_{ms}^v \times \mathcal{M}_{ms}^v|$ represents the influence of distracting-object-related visual tokens.

5 EXPERIMENTS

We evaluate our method across various MLLMs and datasets, comparing it with SOTA hallucination mitigation methods. If not specified otherwise, the experiments default to using greedy sampling.

378 5.1 EXPERIMENTAL SETTING 379

380 Models. We apply and evaluate our methods on widely-used models including LLaVA v1.5-7b, 381 LlaVA v1.6-7b Liu et al. (2024a), InstructBLIP Dai et al. (2023), and mPLUG-Owl2 Ye et al. (2024).

382 **Implementation Details.** The maximum output length is set to 256 across all models, with other 383 model parameters kept at their defaults. To avoid excessive modification, we limit α less than 5 for 384 discriminative tasks and 3 for generative tasks. The early stopping thresholds used in the experiments 385 are as follows: LLaVA-v1.5 and LLaVA-v1.6: 7%, InstrucBLIP: 25%, and mPLUG-Owl2: 2.5%. 386 All experiments are performed on an NVIDIA A40 GPU with batch size of 1.

387 Datasets and Evaluation Metrics. In alignment with established evaluation standards from previ-388 ous studies Zhao et al. (2024); Yin et al. (2023); Zhou et al. (2023); Leng et al. (2023), we evalu-389 ate our hallucination elimination method on discriminative benchmarks, including POPE Li et al. 390 (2023) and the discriminative component of Amber Wang et al. (2023), as well as on open-ended 391 benchmarks, such as the VQA benchmark LLaVa-QA90 Liu et al. (2024b) and image captioning 392 benchmarks from MSCOCO Lin et al. (2014) and Amber Wang et al. (2023). For discriminative tasks, hallucination manifests as a misclassification problem $p(\mathbf{y} = \mathbf{yes} | \mathcal{F}_{\theta}, \mathbf{t}^p \notin \mathbf{t}^v, \mathbf{t}^p)$, we re-393 port both accuracy and F1 score. For open-ended VQA, GPT-4V OpenAI et al. (2024) is used to 394 score both accuracy (Acc) and detailedness (Det) on a scale of 10. For image captioning, we fo-395 cus on object hallucination and report the Caption Hallucination Assessment with Image Relevance 396 (CHAIR) Rohrbach et al. (2018) score, which includes sentence-level (C_S) and instance-level (C_I) 397 percentages, instance-level recall R, and the average generated length (Len). 398

399 5.2 EXPERIMENTAL RESULTS 400

401 Results on POPE. Tab. 1 and Tab.4 in the Appendix show that our methods significantly enhance 402 baseline MLLMs, and outperform SOTA techniques, including the hallucination correction method 403 Woodpecker Yin et al. (2023), and contrastive decoding approaches like VCD Leng et al. (2023), M3ID Favero et al. (2024) and AVISC Woo et al. (2024). These results suggest that explicitly 404 405 and precisely balancing key biases provides substantial benefits. Our improvements are achieved solely through self-reflection and the application of GACD-VA, as GACD-CR and early stopping 406 are not applicable due to the one-word output constraint of the discriminative task. Notably, these 407 improvements are more pronounced on the mPLUG-Owl2 baseline compared to InstructBLIP and 408 LLaVA-v1.5. This disparity may be because InstructBLIP already exhibits a relatively high visual 409 influence ratio, as shown in Fig. 5a. In contrast, LLaVA-v1.5 tends to be overconfident about ob-410 ject existence, having 100% confidence even when the visual influence ratio is less than 30%, as 411 illustrated in Fig. 5b. On the other hand, mPLUG-Owl2 is less overconfident and has a lower visual 412 influence, which likely explains why it shows the greatest improvement. 413

Table 1: Discriminative VQA Comparison on POPE Li et al. (2023) in MSCOCO Adversarial Set-414 ting; Other Results from Zhao et al. (2024), Woo et al. (2024). 415

Method	LL	aVA-v1.5	Ins	tructBLIP	mPl	mPLUG-Owl2		
method	$Acc \uparrow$	F1 ↑	Acc ↑	F1 ↑	Acc ↑	F1 ↑		
Original	79.0	81.1	71.6	74.7	71.5	76.6		
Greedy	79.4	81.6	79.8	81.4	72.5	77.5		
M3ID	77.7	79.7	76.0	77.8	-	-		
AVISC	77.5	79.6	81.6	81.9	-	-		
Woodpecker	80.5	80.6	79.0	78.6	77.5	76.9		
VCD	80.9	81.3	79.6	79.5	-	-		
Ours	83.5	82.1	82.5	82.1	84.2	83.7		

425 **Results on LLaVA-QA90.** We further evaluate our full method, in open-ended VQA tasks on the 426 LLaVA-QA90 Liu et al. (2024b) dataset. The results in Tab. 2 demonstrate the superiority of our 427 approach, consistently improving baseline accuracy and detailness, with at least a 1.44 accuracy 428 improvement and a 0.7 detailness improvement. As illustrated in Fig. 2, in open-ended VQA, our 429 method explicitly amplifies the visual influence to align with the question prompt at the beginning, ensuring that responses are visually grounded. Additionally, it matches the influence of previous 430 outputs as the output length increases, preventing visual forgetting. This explicit balance allows our 431 approach to surpass the VCD method Leng et al. (2023), which also employs contrastive decoding.

Method	LL	aVA-v1.5	In	tructBLIP	mPLUG-Owl2		
litetitot	$Acc \uparrow$	Det ↑	$Acc \uparrow$	Det ↑	$Acc \uparrow$	Det ↑	
Baseline	3.23	3.54	3.84	4.07	4.07	4.33	
VCD	4.15	3.85	4.23	4.69	-	-	
Ours	6.20	5.13	6.28	4.77	6.69	6.28	

Table 2: Open-ended VQA Comparison on LLaVA-QA90; Other Results from Leng et al. (2023).

CHAIR Evaluations on MSCOCO. For the image captioning task, we evaluate our approach on a subset of MSCOCO, following the evaluation settings in Deng et al. (2024). As shown in Tab. 3, our approach consistently and significantly reduces hallucination rates at both the object level (C_S) and image level (C_I) , with particularly strong reductions at the image level (C_I) . This demonstrates its effectiveness in mitigating all types of hallucinations. Hallucination rates are strongly correlated with caption length, where increased length often leads to higher hallucination risks. Compared to SOTA methods, our approach more effectively reduces hallucination rates while maintaining similar caption lengths, demonstrating robustness in balancing detail retention with accuracy. It ensures precise and reliable image captions without compromising between length and hallucination risk.

Method	LLaVA-v1.5			InstructBLIP			mPLUG-Owl2		
incuiou	$C_S\downarrow$	$C_I\downarrow$	$Len\uparrow$	$C_S\downarrow$	$C_I\downarrow$	$Len\uparrow$	$\overline{C_S}\downarrow$	$C_I\downarrow$	$Len\uparrow$
Baseline	48.8	13.4	99.8	57.8	16.5	101.3	59.2	17.6	105.3
OPERA	49.5	13.7	85.7	51.5	15.6	85.8	48.5	16.1	86.1
VCD	44.6	12.5	85.7	63.2	19.5	92.5	51.4	16.0	89.6
Ours	41.0	10.9	85.0	47.4	13.4	93.9	45.0	12.4	83.5

Table 3: Image Captioning Comparison on MSCOCO subset. Other Results from Deng et al. (2024)

Results on Amber. Evaluation on the Amber dataset focuses on both image captioning and comprehensive discriminative performance, covering not only object existence but also categories like attributes and relationships. Tab. 4 and Fig. 6 demonstrate that our method consistently improves the performance of various MLLMs across both tasks. Notably, LLaVA-v1.5 and mPLUG-Owl2 exhibit significant improvements compared to InstructBLIP and LLaVA-v1.6. Our method even enhances LLaVA-v1.5's overall score to outperform LLaVA-v1.6, likely due to differences in the models' original visual influence ratios, as InstructBLIP and LLaVA-v1.6 already have more balanced visual contributions. Additionally, LLaVA-v1.6 emphasizes recall over precision, as indicated by its high scores in both coverage (cov) and recall (R). In terms of discriminative performance, our method achieves a better balance between precision and recall, resulting in improved F1 scores. Furthermore, Fig. 6 shows that our method enhances performance across all categories, with particularly significant gains in existence, attributes, and state. This can be attributed to the increased influence of visual tokens, benefiting categories that are easily discernible from the visual inputs.

Table 4: Results on the AMBER Dataset. Results marked with * are reported in Wang et al. (2023).

	Method	w/Ours	Generative Task				Discriminative Task				Score↑
	inculou	W/Ours	cha↓	$\cot \uparrow$	hal \downarrow	$\cos\downarrow$	acc \uparrow	$P\uparrow$	$R\uparrow$	F1 ↑	Scole
-	LLaVA-v1.5	X* ✓	7.8 6.6	51.0 54.7	36.4 33.0	4.2 2.9	72.0 80.3	93.2 82.9	62.4 89.3	74.7 86.0	83.5 89.7
	LLaVA-v1.6	× ✓	9.9 8.7	56.7 58.3	47.4 43.8	4.3 2.5	80.3 81.2	82.9 85.2	89.3 88.8	86.0 87.0	88.5 89.2
	IntructBLIP	X* ✓	8.8 7.2	52.2 52.1	38.2 34.3	4.4 3.6	76.5 78.1	84.5 88.8	79.0 76.6	81.7 82.2	86.5 87.5
	mPLUG-Owl2	X* ✓	10.6 7.5	52.0 53.6	39.9 34.7	4.5 4.0	75.6 82.1	95.0 87.0	66.9 86.2	78.5 86.6	84.0 89.6

Component analysis in Tab. 5 shows each proposed component contributes to the overall performance. GACD-VA significantly reduces hallucinations while improving object recall. GACD-CR further mitigates co-occurrence bias, a residual form of the text-visual bias addressed by GACD-VA, leading to additional hallucination reduction. The greater reduction at the image level (C_I) compared to the object level (C_S) confirms that addressing this residual bias leads to a more comprehensive



Figure 5: (a) Visual influence ratios across the POPE dataset: the ratios are less than 50% with original MLLMs, GACD successfully increases them to nearly 50%. (b) α value visualization: α inversely correlates with the visual influence ratio, as expected. It drops to 0 when the ratio exceeds 50%, indicating no further adjustments are needed, and also when the ratio is below 50% in cases of 100% prediction confidence.



elimination of all types of hallucinations, improving overall image accuracy. The early-stop mecha nism effectively reduces hallucinations, with only a minor trade-off in object recall.

Fig. 5a shows that GACD-VA successfully increases the overall visual influence to match that of 504 the object-existent question prompt in the POPE dataset ⁴. The visualization of α in Fig. 5b indi-505 cates that the α range of LLaVA-v1.5, InstructBlip, and mPLUG-Owl2 is primarily between 0 and 506 1.5. Furthermore, LLaVA-v1.5 and mPLUG-Owl2 tend to exhibit overconfidence regarding object 507 existence in VQA tasks, as α is set to zero even when the visual influence ratio is below 50%. LLaVA-v1.5 is more prone to overconfidence, with many samples falling within a visual influence 508 range of 15% to 35%, compared to mPLUG-Owl2's range of 20% to 45%. Fig. 2 and Fig.1 in the 509 Appendix show that GACD-VA effectively increases the influence of visual tokens, aligning them 510 with the prompt and previous outputs in open-ended generation tasks. This results in higher predic-511 tion confidence and a reduction in hallucinations. Fig. 3c illustrates the effectiveness of GACD-CR, 512 reducing the hallucination rate of predicting both 'table' or 'chair' in single-object images. 513

Table 5:	Component A	Analysis	Using the	CHAIR Metric
----------	-------------	----------	-----------	--------------

С	ompon	ents		LLaV/	A-v1.5			Instruc	tBLIP		mPLUC	G-Owl2	
VA	CR	ES	$\overline{C_S\downarrow}$	$C_I\downarrow$	$R\uparrow$	$Len \uparrow$	$C_S \downarrow$	$C_I\downarrow$	$R\uparrow$	$Len \uparrow \overline{C_S}$	$s \downarrow C_I \downarrow$	$R\uparrow$	Len
			48.8	13.4	78.6	99.8	57.8	16.5	73.6	101.3 59	.2 17.6	75.8	105.3
1			46.4	11.6	79.0	95.6	53.6	15.1	75.3	108.4 52	.6 14.4	78.2	95.6
1	1		46.2	11.3	79.4	95.5	53.2	14.0	74.6	105.7 52	.3 14.2	78.0	95.5
✓	1	1	41.0	10.9	77.3	85.0	47.4	13.4	72.3	93.9 45	.0 12.4	74.9	83.5

6 CONCLUSION

In this work, we address the challenge of hallucinations in MLLMs by targeting three key biases: text-visual, co-occurrence, and long-term biases. We introduce the GACD framework, which uses gradient-based self-reflection to jointly analyze and balance these biases, mitigating hallucinations without requiring costly resources or tuning. The proposed GACD is the first method capable of fully addressing co-occurrence hallucinations, even when distracting elements are present in the visual input. It amplifies essential visual influences to balance text-visual and co-occurrence biases, and introduces a stopping criterion to mitigate long-term hallucinations.

Our method is limited to white-box MLLMs, as it relies on access to gradients. Its effectiveness
 depends on the importance and clarity of visual information relative to the prompt. Questions about
 object existence are straightforward and primarily rely on visual inputs, whereas questions about
 relationships are less direct and require inference beyond visual inputs. As a post-processing technique, our method does not involve model training. In future work, we aim to explore how insights
 from GACD can guide and improve training strategies for enhanced visual perception in MLLMs.

....

523

524

526

527

528

529

530

531

 ⁴Note that the standard deviation of the visual ratio for InstructBlip, after applying GACD-VA, is small.
 This occurs because InstructBlip is not overly confident, enabling our method to adjust the visual ratio to 50% for nearly all samples.

540 REFERENCES

548

554

574

575

576

- Assaf Ben-Kish, Moran Yanuka, Morris Alper, Raja Giryes, and Hadar Averbuch-Elor.
 Mocha: Multi-objective reinforcement mitigating caption hallucinations. *arXiv preprint arXiv:2312.03631*, 2023.
- Zhiyang Chen, Yousong Zhu, Yufei Zhan, Zhaowen Li, Chaoyang Zhao, Jinqiao Wang, and Ming
 Tang. Mitigating hallucination in visual language models with visual supervision, 2023. URL
 https://arxiv.org/abs/2311.16479.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.
- Ailin Deng, Zhirui Chen, and Bryan Hooi. Seeing is believing: Mitigating hallucination in large vision-language models via clip-guided decoding. *arXiv preprint arXiv:2402.15300*, 2024.
- Alessandro Favero, Luca Zancato, Matthew Trager, Siddharth Choudhary, Pramuditha Perera, Alessandro Achille, Ashwin Swaminathan, and Stefano Soatto. Multi-modal hallucination control by visual information grounding. *arXiv preprint arXiv:2403.14003*, 2024.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu
 Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation
 benchmark for multimodal large language models, 2024.
- 561
 562
 563
 564
 565
 564
 565
 565
 565
 564
 565
 565
 565
 565
 564
 565
 565
 565
 564
 565
 565
 565
 564
 565
 565
 565
 564
 565
 565
 565
 564
 565
 565
 565
 565
 564
 565
 565
 565
 564
 565
 565
 565
 565
 565
 565
 565
 565
 566
 565
 566
 566
 566
 567
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
 568
- Cheongwoong Kang and Jaesik Choi. Impact of co-occurrence on factual knowledge of large lan guage models, 2023.
- Younghyun Kim, Sangwoo Mo, Minkyu Kim, Kyungmin Lee, Jaeho Lee, and Jinwoo Shin. Discovering and mitigating visual biases through keyword explanation, 2024.
- Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing.
 Mitigating object hallucinations in large vision-language models through visual contrastive de coding. arXiv preprint arXiv:2311.16922, 2023.
 - Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.
 Llava-next: Improved reasoning, ocr, and world knowledge, January 2024a. URL https://llava-vl.github.io/blog/2024-01-30-llava-next/.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances in neural information processing systems, 36, 2024b.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey

594 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, 595 Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila 596 Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, 597 Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gib-598 son, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan 600 Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, 601 Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun 602 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-603 mali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook 604 Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel 605 Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen 606 Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel 607 Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, 608 Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, 609 Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, 610 Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel 611 Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Ra-612 jeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, 613 Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel 614 Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe 615 de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, 616 Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, 617 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra 618 Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, 619 Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Sel-620 sam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, 621 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, 622 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-623 ston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vi-624 jayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan 625 Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt 626 Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, 627 Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wo-628 jciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, 629 Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024.

630 631

632

633

634

635

638

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.

- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object
 hallucination in image captioning. *arXiv preprint arXiv:1809.02156*, 2018.
- Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan,
 Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with
 factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023.
- SM Tonmoy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das.
 A comprehensive survey of hallucination mitigation techniques in large language models. *arXiv* preprint arXiv:2401.01313, 2024.
- 646

642

647 Chaojun Wang and Rico Sennrich. On exposure bias, hallucination and domain shift in neural machine translation. *arXiv preprint arXiv:2005.03642*, 2020.

648 649 650	Junyang Wang, Yuhang Wang, Guohai Xu, Jing Zhang, Yukai Gu, Haitao Jia, Ming Yan, Ji Zhang, and Jitao Sang. An llm-free multi-dimensional benchmark for mllms hallucination evaluation. <i>arXiv preprint arXiv:2311.07397</i> , 2023.
652 653 654	Sangmin Woo, Donguk Kim, Jaehyuk Jang, Yubin Choi, and Changick Kim. Don't miss the forest for the trees: Attentional vision calibration for large vision language models, 2024. URL https://arxiv.org/abs/2405.17820.
655 656 657	Shangyu Xing, Fei Zhao, Zhen Wu, Tuo An, Weihao Chen, Chunhui Li, Jianbing Zhang, and Xinyu Dai. Efuf: Efficient fine-grained unlearning framework for mitigating hallucinations in multi-modal large language models. <i>arXiv preprint arXiv:2402.09801</i> , 2024.
658 659 660 661 662	Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, and Fei Huang. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 13040–13051, 2024.
663 664 665	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> , 2023.
666 667	Zihao Yue, Liang Zhang, and Qin Jin. Less is more: Mitigating multimodal hallucination from an eos decision perspective. <i>arXiv preprint arXiv:2402.14545</i> , 2024.
669 670	Bohan Zhai, Shijia Yang, Chenfeng Xu, Sheng Shen, Kurt Keutzer, Chunyuan Li, and Manling Li. Halle-control: Controlling object hallucination in large multimodal models, 2024.
671 672	Linxi Zhao, Yihe Deng, Weitong Zhang, and Quanquan Gu. Mitigating object hallucination in large vision-language models via classifier-free guidance. <i>arXiv preprint arXiv:2402.08680</i> , 2024.
673 674 675 676	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. <i>arXiv preprint arXiv:2310.00754</i> , 2023.
678 679	
680	
681	
682	
683	
684	
685	
697	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	
701	