MITIGATING OBJECT HALLUCINATION IN LARGE VISION-LANGUAGE MODELS VIA IMAGE-GROUNDED GUIDANCE

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028 029

031

Paper under double-blind review

Abstract

The advancement of Large Vision-Language Models (LVLMs) has increasingly highlighted the critical issue of their tendency to hallucinate non-existing objects in the images. To address this issue, previous works focused on using specially curated datasets or powerful LLMs (e.g., GPT-3.5) to rectify the outputs of LVLMs. However, these approaches require either expensive training/fine-tuning or API access to advanced LLMs for post-generation correction. In response to these limitations, we propose Mitigating hallucinAtion via image-gRounded guIdaNcE (MARINE), a framework that is both training-free and API-free. MARINE effectively and efficiently reduces object hallucinations during inference by introducing image-grounded guidance to LVLMs. This is achieved by leveraging open-source vision models to extract object-level information, thereby enhancing the precision of LVLM-generated content. Our framework's flexibility further allows for the integration of multiple vision models, enabling more reliable and robust objectlevel guidance. Through comprehensive evaluations across 5 popular LVLMs with diverse evaluation metrics and benchmarks, we demonstrate the effectiveness of MARINE, which even outperforms existing fine-tuning-based methods. Remarkably, it reduces hallucinations consistently in GPT-4V-assisted evaluation while maintaining the detailedness of LVLMs' generations.

1 INTRODUCTION

The advent of Large Language Models (LLMs) has motivated advancements in extending their remarkable capabilities to multimodal data. Grounded in the development of pre-trained vision-033 language models (Radford et al., 2021; Jia et al., 2021; Alayrac et al., 2022) that align visual and 034 textual embedding spaces, Large Vision Language Models (LVLMs) have gained substantial attention in both architectural development (Liu et al., 2023d; Zhu et al., 2023; Ye et al., 2023; Dai et al., 2023a; Gao et al., 2023), alignment (Yu et al., 2024; Zhou et al., 2024; Deng et al., 2024) and benchmarking 037 datasets (Xu et al., 2023; Lu et al., 2024; Zhang et al., 2024a). However, similar to the hallucination issues in textual LLMs (Ji et al., 2023), where irrelevant content is generated with input prompts, LVLMs face a specific challenge known as object hallucination: generating non-existing objects for a 039 given image (Li et al., 2023b; Wang et al., 2023b; Zhou et al., 2023; Fu et al., 2023; Lovenia et al., 040 2023; Jing et al., 2023). Such a problem is particularly concerning as it compromises the model's 041 accuracy and reliability, especially considering the growing application of LVLMs to safety-critical 042 downstream tasks such as medical imaging (Chambon et al., 2022; Bazi et al., 2023). 043

In response to the pressing issue of object hallucinations in LVLMs, early attempts (Liu et al., 2023a;b; 044 Gunjal et al., 2023; Wang et al., 2023a) focused on addressing the bias by curating high-quality datasets for fine-tuning or leveraging advanced GPT queries (Yin et al., 2023), such as GPT-4, to post-046 process the generated captions. However, these methods can be infeasible to implement. For instance, 047 creating extensive, high-quality datasets for fine-tuning LVLMs is costly and requires significant 048 human annotation. Additionally, relying on advanced GPT models for post-processing is expensive and can raise privacy concerns, especially in sensitive fields like medical imaging. Most importantly, these approaches do not address the *intrinsic* causes of object hallucination in LVLMs. Specifically, 051 fine-tuning simply provides more data for the LVLM to learn, which can lead to overfitting to a particular dataset, as seen with methods like LURE (Zhou et al., 2023). Post-processing methods may 052 also introduce new hallucinations, as they do not inherently correct the root cause of hallucinations in LLMs or LVLMs but just overwrite the generated response.



Figure 1: Illustration of MARINE framework, which introduces a vision toolbox with one or multiple guidance models to enrich the visual context of the original LVLM. The output logits are controlled to place more importance on the guided generation with the guidance strength γ .

In this paper, we investigate the intrinsic causes of object hallucination in LVLMs. Specifically, 071 these deficiencies may stem from the three main components of the LVLMS: 1) insufficient visual 072 context provided by the visual encoder (Zhang et al., 2023b), 2) misalignment between the vision 073 and text domains, and 3) inherent hallucinations common in general language models. To address 074 the first two LVLM-specific causes, we introduce Mitigating hallucinAtion via image-gRounded 075 guIdaNcE (MARINE). MARINE mitigates hallucination issues arising from the visual encoder and 076 domain misalignment by leveraging external guidance from image-grounded models, such as object 077 detection models. Our approach leverages the inherent advantage of these image-grounded models, which are specifically designed and trained for more detailed visual information extraction. These models provide higher quality, fine-grained visual encoding compared to the standard visual encoders 079 in LVLMs, which are primarily optimized for grasping the overall context of an image. Furthermore, we integrate the guidance from image-grounded models into text descriptions, allowing the LVLM to 081 process the information without requiring additional alignment procedures. As a result, MARINE is a training-free, API-free¹ method that addresses object hallucination at inference time by targeting its 083 two root causes. 084

As shown in Figure 1, MARINE incorporates one or more image-grounding models to enrich the 085 visual context of LVLMs. The guidance are then aggregated as prompt input to the LLM decoder to improve the response quality. Empirical evaluations are conducted on five widely-recognized 087 LVLMs across benchmarks including MSCOCO (Lin et al., 2014), LLaVA-QA90 task (Liu et al., 880 2023d), A-OKVOA (Schwenk et al., 2022), and GOA (Hudson & Manning, 2019). We present results based on guidance from a aggregated source of DEtection TRansformer (DETR) (Carion et al., 090 2020) and RAM++ (Huang et al., 2023b). We also include ideal results based on ground truth object oracle, denoted as MARINE-Truth. Our experimental results demonstrate that, in comparison with 091 state-of-the-art algorithms, MARINE exhibits further reduced hallucination, as measured by popular 092 hallucination metrics such as CHAIR (Rohrbach et al., 2018) and POPE (Li et al., 2023b), as well as additional metrics considered in this study including the recall and GPT-4V's evaluation of the 094 responses. These results confirm that MARINE can effectively mitigate object hallucinations without 095 requiring additional training resources or access to advanced LLMs. To summarize, our contribution 096 are listed as follows:

- We introduce MARINE, a universal framework and aggregating a toolbox of image-grounded visual 098 models to guide the generation process of LVLMs. MARINE leverages the intrinsic advantages of these visual models in providing the detailed information of the input image and help mitigate the 100 hallucinations in LVLMs. 101
- Through extensive evaluations on various datasets, we demonstrate that MARINE consistently outperform the baselines in hallucination mitigation while maintaining overall performance across 103 multiple tasks (image captioning, VQA).
 - MARINE provides a favorable trade-off between latency and accuracy, with the lowest computational overhead compared to existing baselines. The minimal increase in latency comparing to the
- 105 107

102

104

067

068

¹The term "API-free" in denotes the elimination of any need for API calls to OpenAI. We note that Woodpecker requires 3-5k input tokens for an API call to each short captioning task.

baselines, combined with the high accuracy of our results, positions MARINE as a practical and
 scalable solution for real-world applications without significant computational cost.

2 RELATED WORK

111

112 **Object Hallucination in Large Vision-Language Models.** Since the introduction of recent Large 113 Vision-Language Models (LVLMs) (Liu et al., 2023d; Zhu et al., 2023; Ye et al., 2023; Dai et al., 114 2023a; Gao et al., 2023), the hallucination phenomenon in these models has gathered significant 115 attention in the research community. This issue was first highlighted by Li et al. (2023b) with 116 subsequent studies (Wang et al., 2023b; Zhou et al., 2023; Fu et al., 2023; Lovenia et al., 2023) 117 that, LVLMs exhibit similar hallucination problems as the textual LLMs. Notably, different from 118 textual LLMs, LVLMs are prone to a unique type of hallucination called 'object hallucination' 119 (Rohrbach et al., 2018), where the model falsely perceives the presence of non-existent objects in 120 images. In response to object hallucination problems, efforts have been made to mitigate object hallucination in smaller image captioning models (Biten et al., 2022; Dai et al., 2023b). Regarding 121 the recent development of LVLMs, several works (Liu et al., 2023b; Gunjal et al., 2023) proposed 122 vision-language fine-tuning datasets aimed for improved robustness. Wang et al. (2023a) leveraged 123 the vision-language model to generate more diverse instruction-tuning data and iteratively correct the 124 inaccuracies in data. Zhai et al. (2023) introduced a GPT-4 assisted evaluation method and also a fine-125 tuning strategy using the MSCOCO dataset. Most related to our setting, Yin et al. (2023) proposed 126 Woodepecker, a five-stage training-free method eventually leveraging GPT-3.5 API for hallucination 127 correction. Concurrently, several works (Leng et al., 2023; Huang et al., 2023a; Chen et al., 2024; 128 Liu et al., 2024; Wan et al., 2024; Zhang et al., 2024b) began to focus on on the training-free setting 129 and can similarly be formulated as approaches using Classifier-Free Guidance (CFG). However, these 130 approaches and MARINE differ in the focus of their method designs within the larger framework of 131 CFG. With detailed discussion deferred to Appendix A, we highlight that MARINE operates at the image level and uniquely employs a vision toolbox to ensemble information from multiple vision 132 models, producing guidance that reaches a consensus among the models for more accurate results. 133 These approaches are further complementary to MARINE, and integrating them could pave the way 134 for a more effective strategy in future research. 135

Controllable Generation. Controllable text generation (Prabhumoye et al., 2020; Hu & Li, 2021; 136 Zhang et al., 2023a) has emerged as a vital research domain, focusing on the generation of natural 137 sentences with controllable attributes such as persona (Prabhumoye et al., 2020; Hu & Li, 2021; 138 Zhang et al., 2023a)and politeness (Niu & Bansal, 2018; Madaan et al., 2020). Among the various 139 approaches, fine-tuning has been recognized as the most straightforward approach, achieved either 140 through full fine-tuning (Li & Liang, 2021; Ouyang et al., 2022; Carlsson et al., 2022) or integrating 141 tunable adaptors (Lin et al., 2021; Ribeiro et al., 2021). While fine-tuning has been effective in 142 a wide range of applications, it is also expensive in computation as the size of LLMs is growing 143 tremendously. Recently, there has been a development on controllable generation with diffusion 144 models (Li et al., 2022; Lin et al., 2023b), extending to controllable text-to-image generation (Yang 145 et al., 2023). Particularly, the use of classifier guidance (Dhariwal & Nichol, 2021) and classifier-free guidance (Ho & Salimans, 2021) has become prominent in refining the quality of generated outputs. 146 Most recently, Sanchez et al. (2023) applied classifier-free guidance to language models in the 147 single-modal setting to improve their performance at inference time. Our approach methodologically 148 resembles classifier-free guidance for LVLMs' text generation, while specifically addressing the 149 *multi-modal* context and focusing on reducing hallucinations. 150

3 PRELIMINARIES

151

152

Notation. We use lower case letters, lower case bold face letters, and upper case bold face letters to denote scalars, vectors, and matrices respectively. We use the symbol p to represent the conditional probability of LLM's response. And we denote the sequence of tokens generated before the t-th token as $\mathbf{y}_{< t} = [y_1, \dots, y_{t-1}]$ for t > 1. $\mathbf{y}_{< t}$ is an empty sequence when t = 1.

157 Generative language models. Let p_{θ} denotes an LLM parameterized by θ . Consider a sequence **158** $\mathbf{x} = [x_1, \dots, x_n]$ as the input prompt, where each x_i is a token from a predefined vocabulary. The **159** LLM then generates the response sequence $\mathbf{y} = [y_1, \dots, y_m]$ by sampling from the conditional **160** probability distribution $p_{\theta}(\cdot|\mathbf{x})$, where y_t denotes individual token for $1 \le t \le m$. The conditional **161** distribution $p_{\theta}(\mathbf{y}|\mathbf{x})$ can therefore be expressed as $p_{\theta}(\mathbf{y}|\mathbf{x}) = \prod_{t=1}^m p_{\theta}(y_t|\mathbf{x}, \mathbf{y}_{<t})$, where $\mathbf{y}_{<t} = [y_1, \dots, y_{t-1}]$ for t > 1 and is empty for t = 1. In the case of LVLMs, visual tokens $\mathbf{b} = [v_1, \dots, v_k]$ are additionally included. These tokens are generated from a pre-trained visual encoder and mapped into the token space through a linear projection. The conditional distribution of output y given the visual tokens b and textual prompt x is expressed as $p_{\theta}(\mathbf{y}|\mathbf{b}, \mathbf{x}) = \prod_{t=1}^{m} p_{\theta}(y_t|\mathbf{b}, \mathbf{x}, \mathbf{y}_{< t})$, where p_{θ} is approximated by LVLMs.

Guidance in generative models. The process of a guided generation involves getting the output y conditioned on input x, which encodes the desired properties of the output y. This guidance can be generally added to the model by two distinct approaches: classifier guidance (Dhariwal & Nichol, 2021) and classifier-free guidance (Ho & Salimans, 2021). As a top-level view, both methods formulate the conditional probability distribution of output y conditioned on guidance x as

$$p(\mathbf{y}|\mathbf{x}) \propto p_{\boldsymbol{\theta}}(\mathbf{y}) p(\mathbf{x}|\mathbf{y})^{\gamma},$$
 (3.1)

173 where $p_{\theta}(\mathbf{y})$ is the original generative model and $p(\mathbf{x}|\mathbf{y})$ is the posterior distribution of \mathbf{x} given \mathbf{y} and 174 γ is the guidance strength. In the classifier guidance, the posterior distribution $p(\mathbf{x}|\mathbf{y})$ in equation 3.1 175 is replaced by a classifier $p_{\phi}(\mathbf{x}|\mathbf{y})$ parameterized by ϕ , which requires additional training step and 176 calculating $\nabla_{\mathbf{x}} \log p_{\boldsymbol{\phi}}(\mathbf{x}|\mathbf{y})$. The classifier-free guidance, on the other hand, removes the necessity of 177 the parameterized classifier f_{ϕ} . Instead, according to the Bayes rule, the posterior distribution can be 178 approximated by $p_{\theta}(\mathbf{x}|\mathbf{y}) \propto p_{\theta}(\mathbf{y}|\mathbf{x})/p_{\theta}(\mathbf{y})$, where $p_{\theta}(\mathbf{y}|\mathbf{x})$ is the generative model when taking 179 \mathbf{x} as prompt input. Plugging this back into equation 3.1 yields the guided distribution that can be 180 approximated by

181 182

187

$$\widehat{p}_{\theta}(\mathbf{y}|\mathbf{x}) \propto p_{\theta}(\mathbf{y}) \cdot p_{\theta}(\mathbf{y}|\mathbf{x})^{\gamma} / p_{\theta}(\mathbf{y})^{\gamma} = p_{\theta}(\mathbf{y}|\mathbf{x})^{\gamma} / p_{\theta}(\mathbf{y})^{\gamma-1}.$$

As a result, the guided LLM \hat{p}_{θ} places more importance on the prompt x during generation with the increasing value of γ , thereby producing texts that better align with the desired behavior from the prompt (Sanchez et al., 2023).

4 Method

188 The existing architecture of LVLMs is usually composed of a visual encoder, a visual and textual 189 domain alignment layer, and the LLM itself. Therefore, besides the inherent language priors of 190 LLMs (Biten et al., 2022), object hallucination may arise from (1) deficiencies in the visual encoder 191 providing insufficient visual information (Zhang et al., 2023b) and (2) misalignment between the 192 visual and textual domains. To mitigate object hallucinations, we introduce MARINE, a framework 193 containing two major components to address the aforementioned challenges: (1) introducing additional visual information from a set of vision models and (2) using the additional aggregated visual 194 features to guide the LVLM's generation. In Figure 1, we present the framework overview. 195

196 197

4.1 VISUAL GUIDANCE FROM IMAGE-GROUNDED FEATURES

To introduce image-grounded guidance to mitigate hallucinations, our approach integrates additional 198 object detection models, which differ from the visual encoders used in LVLM that are usually pre-199 trained from CLIP (Radford et al., 2021). This integration leverages the object detection models 200 to extract detailed visual information from images. Upon acquiring these extra visual information 201 from different image-grounded models, we aggregate and translate the collected information into 202 textual information. This aggregation can be done by the language model (Lin et al., 2023a) or rule 203 based algorithm (Bird et al., 2009). Such an information aggregation is effective and efficient, as 204 it eliminates the necessity of fine-tuning the alignment layer while retaining the rich information 205 encoded by various of image grounding models. We subsequently employ a simple prompt "focusing 206 on the visible objects in this image:" and concatenate it with the aggregated object information, 207 denoted as the guidance prompt c.

 208
 309
 4.2
 Guided Text Generation with Visual Information

210 We tackle the object hallucination problem of LVLMs by specifically placing importance on the 211 additional image-grounded information we introduced. In addition to the visual tokens *b* extracted 212 from the original LVLM and textual prompt x, we extract the auxiliary visual tokens c from the 213 additional guidance models. The generation of the *t*-th token in the output y of our classifier-free 214 guided LVLM p_{θ} is expressed as

$$\widehat{p}_{\theta}(y_t|\boldsymbol{b}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{< t}) \propto p_{\theta}(y_t|\boldsymbol{b}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{< t})^{\gamma} / p_{\theta}(y_t|\boldsymbol{b}, \mathbf{x}, \mathbf{y}_{< t})^{\gamma-1},$$

where c denotes our control guidance and γ is the control strength. The sampling of output generation is given by

218 219

$$\widehat{p}_{\theta}(\mathbf{y}|\boldsymbol{b}, \mathbf{c}, \mathbf{x}) = \prod_{t=1}^{m} \widehat{p}_{\theta}(y_t|\boldsymbol{b}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{< t}) \propto \prod_{t=1}^{m} \frac{p_{\theta}(y_t|\boldsymbol{b}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{< t})^{\gamma}}{p_{\theta}(y_t|\boldsymbol{b}, \mathbf{x}, \mathbf{y}_{< t})^{\gamma-1}} = \frac{p_{\theta}(\mathbf{y}|\boldsymbol{b}, \mathbf{c}, \mathbf{x})^{\gamma}}{p_{\theta}(\mathbf{y}|\boldsymbol{b}, \mathbf{x})^{\gamma-1}} = \frac{p_{\theta}(\mathbf{y}|\boldsymbol{b}, \mathbf{c}, \mathbf{x})^{\gamma}}{p_{\theta}(\mathbf{y}|\boldsymbol{b}, \mathbf{x})^{\gamma-1}}$$

We can further view MARINE in the logit space, where the *t*-th token is therefore sampled from the logit space by

$$\log \widehat{p}_{\boldsymbol{\theta}}(y_t | \boldsymbol{b}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{< t}) = \gamma \log p_{\boldsymbol{\theta}}(\mathbf{y} | \boldsymbol{b}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{< t}) + (1 - \gamma) \log p_{\boldsymbol{\theta}}(\mathbf{y} | \boldsymbol{b}, \mathbf{x}, \mathbf{y}_{< t}).$$

224 This linear combination of logits implies that the conditional generation on the additional image-225 grounded guidance acts as a controllable gate. Only objects with relatively high probabilities in 226 both branches could appear at top when sampling. Specifically, setting $\gamma = 0$ recovers the original 227 LLM generation without control guidance and setting $\gamma = 1$ produces the LLM generation entirely 228 based on the control. Meanwhile, for $\gamma \in (0, 1)$, MARINE yields a combination of the original 229 generation $p_{\theta}(\mathbf{y}|\mathbf{b}, \mathbf{x})$ and the generation conditioned on the guidance $p_{\theta}(\mathbf{y}|\mathbf{b}, \mathbf{c}, \mathbf{x})$. This strikes 230 a balance between a better ability to follow instructions to generate high-quality answers and the 231 increased accuracy and detail in image descriptions. The formulation therefore shares resemblance to 232 the classifier-free guidance introduced for LLMs (Sanchez et al., 2023), which places importance on 233 the textual prompt itself to better align the LLM generation with user intention in the *single-modal* setting. We summarize MARINE in Algorithm 1. In detail, MARINE aggregates the collected visual 234 information $\{c_i\}_i$ using function Aggr., which can be a small language model for information 235 aggregation (Lin et al., 2023a), or a rule-based algorithms like majority voting (as similarly used by 236 Wang et al.). Notably, MARINE only double the LLM inference time of in Line 7 and Line 9, while 237 adding the guidance from each single image grounded model will significantly increase the inference 238 time when the number of image grounded models increase. 239

Algorithm 1 Mitigating hallucinAtion via image-gRounded guIdaNcE (MARINE)

240 1: Input: LLM parameter θ , input prompt x, visual tokens b from LVLM's original vision tower 241 2: Input: auxiliary visual tokens $\{\mathbf{c}_i\}_{i=1}^M$ from M image grounding models, guidance scale γ 242 3: Initialize empty output $\mathbf{y} = []$. 243 4: Aggregate visual information as textual prompt $\mathbf{c} = \text{Aggr.}(\{\mathbf{c}_i\}_{i=1}^M)$ 244 5: for t = 0, 1, ..., T do 245 Construct unconditional input $\mathbf{x}_{uncond}^{(t)} = [\mathbf{b}, \mathbf{x}, \mathbf{y}_{< t}]$. 6: 246 Generate unconditional output logits using LLM: $\ell_{\text{uncond}}^{(t)} = \log p_{\theta}(\mathbf{x}_{\text{uncond}}^{(t)})$. 7: 247 Construct conditional input $\mathbf{x}_{cond}^{(t)} = [\mathbf{b}, \mathbf{c}, \mathbf{x}, \mathbf{y}_{< t}].$ 8: 248 Generate conditional output logits using LLM: $\ell_{\text{cond}}^{(t)} = \log p_{\theta}(\mathbf{x}_{\text{cond}}^{(t)})$. 249 9: Update output logits $\ell^{(t)} = \gamma \ell^{(t)}_{\text{cond}} + (1 - \gamma) \ell^{(t)}_{\text{uncond}}$. 250 10: 251 Sample token y_t from logit space denoted by $\ell^{(t)}$ 11: 252 12: Let $\mathbf{y} = [\mathbf{y}, y_t]$. 253 13: end for 254 14: **Output:** y.

5 EXPERIMENTS

255 256

257

In this section, we evaluate MARINE in mitigating object hallucinations across various LVLMs,
 showing that it outperforms state-of-the-art methods on established metrics across different question
 formats.

261 5.1 EXPERIMENT SETUP 262

Models. To demonstrate the broad applicability of our approach across different LVLM architectures, we apply and evaluate MARINE to recent widely-used models including *LLaVA* (Liu et al., 2023d), *LLaVA-v1.5* (Liu et al., 2023c), *MiniGPT-v2* (Chen et al., 2023), *mPLUG-Owl2* (Ye et al., 2023) and *InstructBLIP* (Liu et al., 2023c). To address the object hallucination problems in text generation, we incorporate the DEtection TRansformer (DETR) (Carion et al., 2020) and RAM++ (Huang et al., 2023b) as the additional vision models for guidance.

Guidance from Multiple Sources. Our framework's compatibility with various vision models allows for the incorporation of multiple sources to enhance precision and robustness. By considering

270 object-level information from DETR and RAM++ simultaneously, we generate guidance that reflects 271 consensus across these models. This approach significantly improves the accuracy and reliability of 272 the guidance provided to the LVLM. 273

Datasets and evaluations. In alignment with established evaluations from previous studies (Dai et al., 2023b; Yin et al., 2023), we assess our method using the following metrics:

• Caption Hallucination Assessment with Image Relevance (CHAIR) (Rohrbach et al., 2018). It involves prompting the LVLMs to generate a description for the input image, and then comparing this generation with ground truth objects present in the image. CHAIR quantifies hallucination both at instance level and sentence level, respectively defined as $CHAIR_I$ and $CHAIR_S$:

$$CHAIR_{I} = \frac{\left| \{ \text{hallucinated objects} \} \right|}{\left| \{ \text{all mentioned objects} \} \right|}, \quad CHAIR_{S} = \frac{\left| \{ \text{captions with hallucinated objects} \} \right|}{\left| \{ \text{all captions} \} \right|}.$$

In addition to these metrics, we incorporate an instance-level Recall score in our evaluation to evaluate whether the descriptions accurately include the necessary visual content from the image:

Recall = $|\{\text{non-hallucinated objects}\}|/|\{\text{all existing objects}\}|$.

284 285

274

275

276

277

278

279

281 282

283

286 287

288

289

290

291

292

293

294

295

296

297

298

299

302

303

304

305

- Polling-based Object Probing Evaluation (POPE) (Li et al., 2023b). POPE formulates a binary classification task by prompting LVLMs with questions such as "Is there a keyboard in this image?" to answer "yes" or "no". We specifically focus on the adversarial setting, which is considered the most challenging setting. Results for the random and popular settings are detailed in Appendix E. We report the accuracy and F1 score of the LVLMs' responses, and the proportion of "yes" answers.
- GPT-4V-aided Evaluation (Yin et al., 2023). The GPT-4V-aided evaluation compares the outputs of two LVLM assistants using GPT-4V as a judge. In this evaluation, we utilize the LLaVA-QA90 task (Liu et al., 2023d)² (including conversations, visual perceptions, and complex reasoning tasks) and additionally consider the image captioning task.

Consistent with Li et al. (2023b), we randomly sampled a subset of 500 images from MSCOCO (Lin et al., 2014) dataset for CHAIR evaluation. For the POPE evaluation, we created 3000 questions across three datasets—500 images each from MSCOCO, A-OKVQA (Schwenk et al., 2022), and GQA (Hudson & Manning, 2019). For the GPT-4V-aided evaluation, we utilized 90 questions from the LLaVA-QA90 task and randomly selected 50 MSCOCO images for image captioning task.

300 Baselines. In addition to comparing with the performance of the original LVLM sampling method, 301 we also consider the following popular methods for mitigating hallucinations.

- *Greedy-Decoding*, which adopts the greedy sampling strategy, by generating tokens with the highest posterior probability to address hallucinations arising from.
- LURE (Zhou et al., 2023), which identifies and masks potentially hallucinated words and fine-tune a MiniGPT4 model to rectify object hallucinations in the generated descriptions.
- LURE with Cutoff. The original LURE method tends to generate long descriptions regardless of 306 the provided instructions, which sometimes results in even worse performance as unnecessary 307 information is included. Therefore, we also introduce a modified baseline, where we truncate the 308 LURE's output to match the length (in terms of the number of sentences) of the original generations.
- Woodpecker (Yin et al., 2023), which leverages GPT-3.5 to correct hallucinations in LVLM 310 generation with five steps toward the correction. 311
 - VCD (Leng et al., 2023), which distorts the image inputs to impose penalties on logit outputs.
- 312 • OPERA (Huang et al., 2023a), which penalizes logits to mitigate over-trust in beam-search decoding 313 and adjusts token selection.
- 314 Lastly, the performance of MARINE improves in correlation with the advancement of the control 315 guidance extractor used. Consequently, to demonstrate the potential upper bound of MARINE's 316 performance, we consider a version utilizing a ground-truth oracle extractor, which we denote as MARINE-Truth. Further details on model architectures, datasets and evaluation metrics are deferred 317 to Appendix C. 318
- 319 **Hyperparameter Setting.** The hyperparameters for our method are fixed across tasks, with key 320 settings including a guidance strength of 0.7, noise intensity for DETR at 0.95, a detection threshold 321 for RAM++ of 0.68, and a greedy sampling approach with a random seed of 242.
- 322 323

²https://github.com/haotian-liu/LLaVA/blob/main/playground/data/ coco2014_val_gpt4_qa_30x3.jsonl

324 Table 1: Evaluation with CHAIR score across multiple LVLM architectures comparing our method 325 with several baselines. We report CHAIR_S, CHAIR_I and the recall score. The **bold** numbers indicate 326 the best results among the methods evaluated and the underscored numbers represent the second-best 327 results. We show MARINE-Truth as a reference performance of MARINE.

28	Method	1	LLaVA		LL	aVA-v1	.5	М	iniGPTv	2	mPI	LUG-Ov	vl2	Ins	tructBL	IP	A	verage	
9	CHAIR	$C_S \downarrow$	$C_I\downarrow$	$R\uparrow$	$C_S \downarrow$	$C_I\downarrow$	$R\uparrow$	$C_S \downarrow$	$C_I \downarrow$	$R\uparrow$	$C_S \downarrow$	$C_I\downarrow$	$R\uparrow$	$C_S \downarrow$	$C_I\downarrow$	$R\uparrow$	$C_S \downarrow$	$C_I\downarrow$	$R\uparrow$
)	Greedy	26.6	10.5	47.4	8.8	4.6	41.1	8.2	<u>4.2</u>	41.1	6.2	3.4	38.8	5.0	3.2	33.2	11.0	5.2	40.3
	LURE	33.8	11.6	54.8	38.9	11.2	56.3	36.2	11.4	54.6	33.9	10.8	55.9	38.1	12.1	54.5	36.2	11.4	55.2
	LURE w/ cutoff	24.4	9.3	50.2	18.4	6.8	<u>47.3</u>	12.5	6.2	42.0	15.4	6.6	<u>45.5</u>	9.6	6.4	34.5	16.1	7.1	43.9
	Woodpecker	19.5	<u>8.9</u>	44.3	8.5	4.5	38.4	7.5	4.5	37.0	8.0	4.3	37.5	8.0	6.2	32.6	10.3	5.7	38.0
	VCD	28.1	11.0	46.6	<u>7.3</u>	<u>4.1</u>	40.8	6.8	3.9	38.2	<u>5.9</u>	3.4	37.7	<u>2.4</u>	<u>1.5</u>	33.7	<u>10.1</u>	<u>4.8</u>	39.4
	OPERA	<u>22.4</u>	9.9	43.6	11.0	6.7	40.2	9.2	5.0	41.3	5.8	<u>3.2</u>	38.4	4.6	2.7	<u>38.0</u>	10.6	5.5	40.3
	MARINE	17.8	7.2	<u>50.8</u>	6.2	3.0	44.3	11.8	4.9	<u>49.7</u>	4.2	2.3	41.4	2.2	1.3	36.3	8.4	3.7	<u>44.5</u>
)	MARINE-Truth	19.6	5.1	79.0	6.0	2.5	55.3	12.6	3.8	70.5	3.8	1.7	48.0	3.0	1.8	35.9	8.9	2.9	57.5

Table 2: Evaluation with POPE score in adversarial setting across multiple LVLM architectures comparing our method with several baselines. We report the POPE accuracy (%), F1 score (%) and the yes ratio (%). The ideal yes ratio for a non-biased LVLM is 50%. The **bold** numbers indicate the best results among the methods evaluated and the underscored numbers represent the second-best 340 results. We show MARINE-Truth as a reference performance of MARINE.

341									1										
0.10	Method	1	LLaVA		LL	aVA-v1	.5	Mi	niGPTv	2	mPl	LUG-Ow	/12	Inst	ructBL	IP	А	verage	
342	POPE	$\mathrm{Acc}\uparrow$	F1 ↑	Yes	$Acc\uparrow$	F1↑	Yes	$\rm Acc\uparrow$	F1 ↑	Yes	$Acc\uparrow$	$F1\uparrow$	Yes	$\operatorname{Acc}\uparrow$	F1 ↑	Yes	Acc ↑	F1 ↑	Yes
343	Greedy	51.8	67.4	97.7	79.4	<u>81.6</u>	61.6	82.7	81.7	<u>44.5</u>	72.5	77.5	72.4	<u>79.8</u>	81.4	58.6	73.2	77.9	67.0
344	LURE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
245	Woodpecker	77.5	77.6	50.5	80.5	80.6	50.5	79.5	77.8	42.5	77.5	76.9	<u>47.5</u>	79.0	78.6	48.0	<u>78.8</u>	<u>78.3</u>	<u>47.8</u>
340	VCD	54.6	68.5	94.0	78.2	80.7	62.8	81.4	80.2	44.1	72.3	77.0	70.5	79.7	80.9	<u>56.7</u>	73.2	77.5	65.6
346	OPERA	51.7	67.4	98.0	77.5	80.1	63.2	<u>82.9</u>	<u>81.9</u>	44.3	70.3	<u>79.1</u>	84.6	<u>79.8</u>	81.4	58.6	72.4	78.0	69.7
347	MARINE	<u>66.9</u>	<u>72.9</u>	<u>72.3</u>	85.0	84.3	<u>45.7</u>	83.0	82.9	49.4	82.8	82.7	49.2	81.7	79.4	38.8	79.9	80.4	51.1
348	MARINE-Truth	75.6	80.1	72.3	92.0	92.5	57.0	86.9	88.3	62.5	93.4	93.76	56.2	93.8	93.8	51.0	88.3	89.7	59.8
~																			

5.2 RESULTS

337

338

339

349 350

351

352

353

Experimental results on object hallucination metrics (CHAIR and POPE) are presented in Table 1 and Table 2. Overall, MARINE achieves superior performances across different LVLM architectures and evaluation metrics, ranked as the best or second-best on the majority of the evaluation metrics.

Results on CHAIR. Table 1 presents the evaluation of various mitigation methods using CHAIR 354 scores across multiple LVLM architectures. The results demonstrate that MARINE consistently 355 outperforms other state-of-the-art methods, achieving the highest average scores in both $CHAIR_S$ and 356 CHAIR_I and the second-best Recall score. Specifically, MARINE surpasses the second-best perform-357 ing method by an average margin of 1.7 on CHAIR_S and 1.1 on CHAIR_I. Notably, MARINE exhibits 358 exceptional performance on LLaVA architectures, with improvements in CHAIR scores of up to 8.8 359 compared to its original performance. In contrast, methods such as LURE and Woodpecker show less 360 effectiveness in hallucination mitigation. The reference method, MARINE-Truth, generally achieves 361 the strongest results, as expected given its access to ground-truth guidance. However, MARINE's per-362 formance closely approximates that of its ground-truth counterpart, indicating successful leveraging 363 of multiple guidance models to provide reliable control in LVLM generation.

364 **Results on POPE.** The POPE evaluation, presented in Table 2, further validates the superior performance of MARINE against existing baselines across various question formats. MARINE consistently 366 outperforms all other methods by a substantial margin, demonstrating average improvements of 6.7%367 in accuracy and 3.5% in F1 score relative to the original outputs across models. Even when compared 368 to the second-best method, Woodpecker, MARINE maintains a performance edge of 1.1% and 2.1%369 respectively in accuracy and F1 score. Moreover, MARINE effectively mitigates the LVLMs' biased 370 tendency towards affirmative responses, as evidenced by a more balanced "yes" ratio (closer to 50%, representing a 15.9% shift towards unbiased answers). This improvement notably addresses the 371 overconfidence issue prevalent in existing models. 372

373 **Results on GPT-4V-aided evaluation.** Following Yin et al. (2023), we leverage GPT- $4V^3$ to 374 evaluate and compare the performance of the original LVLMs and LVLMs with MARINE on LLaVA-375 QA90 and an image captioning task. This GPT-4V-assisted evaluation introduces a qualitative

³We used gpt-4-1106-vision-preview in obtaining our final experiment results. As OpenAI continues to update its API, different versions may result in slightly different values.

Table 3: Results of GPT-4V-aided evaluation. The accuracy and detailedness metrics are on a scale of 10, and a higher score indicates better performance. The symbols \times and \checkmark indicate performance metrics without and with our method, respectively.

Taala	Matriag	LLa	iVA	mPLUG-Owl2			
Task	Metrics	X	\checkmark	X	\checkmark		
	Acc ↑	$5.82_{\pm 0.10}$	5.94 $_{\pm 0.05}$	$6.03_{\pm 0.13}$	$\textbf{6.35}_{\pm 0.21}$		
LLavA-QA90	Detail ↑	$4.59_{\pm 0.08}$	$4.59_{\pm0.08}$	$5.06_{\pm 0.05}$	$\textbf{5.16}_{\pm 0.10}$		
Imaga Cantioning	Acc \uparrow	$5.27_{\pm 0.20}$	$\textbf{6.11}_{\pm 0.23}$	$7.97_{\pm 0.25}$	$\textbf{8.63}_{\pm 0.20}$		
inage Captioning	Detail ↑	$\textbf{4.39}_{\pm 0.29}$	$4.36_{\pm0.17}$	$5.74_{\pm 0.24}$	$\textbf{6.19}_{\pm 0.23}$		

Table 4: POPE results across three datasets. We report the average score under random, popular, adversarial settings. The detailed POPE results can be found in the appendix **E**. The **bold** numbers indicate the best results. The ideal yes ratio for a non-biased LVLM is 50%.

Datasat	MADINE	L	LaVA		mPLUG-Owl2		
Dataset	W/MARINE	Accuracy ↑	F1↑	Yes(%)	Accuracy ↑	$F1\uparrow$	Yes(%)
MSCOCO (Lin et al. 2014)	X	54.2	68.5	95.5	76.7	80.4	68.2
MSCOCO (Lin et al., 2014)	\checkmark	72.2	76.4	66.9	85.5	85.0	46.5
A OKVOA (Soloverk et al. 2022)	×	51.8	67.5	97.9	69.6	76.5	78.5
A-OK VQA (Schwenk et al., 2022)	\checkmark	64.3	72.8	80.2	82.0	83.5	57.2
COA (Hudson & Manning 2010)	×	52.0	67.6	97.8	73.7	78.7	72.6
GQA (Hudson & Manning, 2019)	\checkmark	62.5	71.8	81.8	80.1	80.6	51.1

perspective beyond the numerical metrics of CHAIR and POPE, offering a richer assessment of
 model performance. The evaluation prompt is detailed in Appendix C.5. As shown in Table 3, GPT-4V
 consistently assigns higher accuracy with equal detailedness scores to models enhanced by MARINE,
 highlighting its ability to produce more precise and detailed descriptions, which demonstrates the
 robustness of our method in real-world visual tasks.

Additional Results on Other Vision-Language Tasks. To further evaluate the generalizability of our approach beyond object hallucination and the MSCOCO dataset, we extended our evaluations to additional datasets including A-OKVQA and GQA and included more general caption quality metrics. As shown in Table 4, the POPE results on datasets such as MSCOCO, A-OKVQA, and GQA demonstrate that our method consistently mitigates hallucinations across various datasets with different image distributions. Figure 2 presents a comprehensive evaluation of the image captioning task on MSCOCO and LLaVA-QA90, a comprehensive VQA dataset, using metrics including BLEU(Papineni et al., 2002), ROUGE(Lin, 2004), CIDEr(Vedantam et al., 2015) and SPICE(Anderson et al., 2016). These results demonstrate that, although our method primarily targets hallucination mitigation, it maintains the overall performance of LVLMs on broader tasks, with no significant trade-offs in caption or VQA quality.



Figure 2: MARINE leads to consistent enhancement in the text qualities on general metrics. Dashed
 lines and solid lines represent without or with MARINE. Higher scores indicate better quality and
 greater similarity between the generated captions and the reference texts.

Table 5: Inference Latency Comparison. We report both the latency and the ratio to the latency of greedy decoding of the original LVLM model.

		Greedy	LU	RE	W	/oodpecl	cer*	VCD	OP	ERA MZ	ARINE (ours)
Training (Cost	0	10min on	A100 8	0G	0		0		0	0
Inference	Latency ^(ms/token)	26.3 (×1.0)	179.9	(×6.84)	9	4.5 (×3.5	59)* 5	3.4 (×2.03) 185.1	(×7.0)	52.2 (×1.98)
*Woodpeck	er requires GPT A	PI key access an	d the latenc	y may d	epend or	OPENA	I API.				
Table 6: using ind reliable a	Ablation stud dividual visio and robust ob	dy comparin on models. ' oject-level g	ig the per This app uidance,	rform roach resul	ance o 1 lever ting ir	f coml ages n 1 super	bining hultipl rior pe	DETR e objec rformar	and RA t detect the on (M++ mo tors to pro CHAIR r	dels versus ovide more netrics.
	I	Model		LLa	aVA	LLaV	A-v1.5	mPLUC	G-Owl2		
		CHAIR		$C_S \downarrow$	$C_I\downarrow$	$C_S \downarrow$	$C_I\downarrow$	$C_S \downarrow$	$C_I\downarrow$		
	1	Ensembling Mo	dels								
	1	MARINE		17.8	7.2	6.2	3.0	4.2	2.3		
	5	Single Models									
	ľ	MARINE-DETH	R only	27.6	8.4	10.5	4.3	5.3	2.7		
	11	MARINE-RAM	only	29.0	9.1	6.6	3.7	5.2	2.8		
	Table	7: Effect of	Integrat	ion M	lethod	s for I	mage-	Ground	ing Mo	odels.	
	Model				LLa	A	LLaV	A-v1.5	mPLU	UG-Owl2	
	CHAIR			\overline{C}	$s \downarrow$	$C_I \downarrow$	$\overline{C_S\downarrow}$	$C_I\downarrow$	$\overline{C_S \downarrow}$	$C_I \downarrow$	
	MARINE-i	Intersect	ion (our	s) 1	7.8	7.2	6.2	3.0	4.2	2.3	-
	MARINE-u	union		3	0.4	9.7	5.4	2.7	4.8	2.7	
											-

Latency Analysis Mitigating object hallucination often requires additional computational resources, 455 a characteristic common to many existing methods which typically involve additional post-generation 456 correction models (Zhou et al., 2023; Zhai et al., 2023; Yin et al., 2023), object detectors (Yin 457 et al., 2023), or more complex decoding processes (Huang et al., 2023a; Leng et al., 2023) to 458 reduce hallucinations. Furthermore, to assess the practical feasibility of our approach in terms 459 of computational costs, we have compared our method with existing baselines on LLaVA-7B. As 460 demonstrated in Table 5, our method increases the decoding time by factors of 1.98, which is the 461 lowest costs among existing baselines, suggesting MARINE can be widely applied with negligible 462 cost. Our method offers the most favorable trade-offs between latency and accuracy in hallucination 463 mitigation. Detailed experiment setting is in Appendix C.6.

465 5.3 ABLATION STUDY

464

How Does Incorporating Multiple Sources to Form Guidance Impact Performance? We
 perform an ablation study to assess the impact of incorporating DETR and RAM++ compared to
 using each model individually, as presented in Table 6. Notably, DETR allows for highly accurate
 object detection, while RAM++ excels in extensive recognition tasks, adding fine-grained details
 to image understanding. Combining the strengths of these image-grounding models, we achieve
 significant performance improvements on the CHAIR metrics. This demonstrates that leveraging
 complementary visual contexts can substantially enhance overall model effectiveness.

Which Method of Integrating Image-Grounding Models Works Best? We investigate two approaches for integrating image-grounding models: using either the intersection or union of detected objects. As shown in Table 7, the intersection-based method outperforms the union, significantly reducing object hallucination. This result highlights the importance of precision and consistency in guidance, as taking intersection ensures consensus across models, leading to more reliable guidance. The detailed experimental setup and prompt templates are provided in Appendix C.

Effect of Guidance Strength. Figure 3 shows that increasing guidance strength from 0 to 1 leads to a notable decrease in CHAIR scores. This trend suggests that higher guidance strength makes LVLMs rely more on image-grounded features provided by DETR, thereby enhancing their ability to produce accurate descriptions. It's crucial to note that, although some models exhibit optimal performance at a guidance strength of $\gamma = 1$, excessively strong guidance can adversely affect the models' ability to adhere to provided instructions. Experimental evidence is detailed in Appendix E. This observation highlights the necessity of having a balanced guidance strength that ensures high-quality, accurate outputs while adhering closely to the given instructions. Based on our findings, we recommend



Figure 3: Ablation study on the effect of guidance strength (γ) on the performance of LLaVA and mPLUG-Owl2 using CHAIR metrics, with γ ranging from 0 to 1.

a guidance strength within the range of $\gamma \in (0.3, 0.7)$ as the most effective for maintaining this balance.

Examples of MARINE's Guided Generation. In Figure 4, we provide specific generation examples
of LLaVA based on queries from different tasks, with or without MARINE. In the first example,
LLaVA incorrectly identifies a white chair in the image, an instance of hallucination as the object
present is a white bird instead. In contrast, MARINE successfully mitigates this hallucination, correctly
guiding the the model to recognize the object as a white bird. Similarly, in the second example,
LLaVA erroneously state that the skateboard rider is holding onto the trucks. With MARINE, the
model's response is more accurate, focusing on verifiable visual elements and correctly stating that
the person is standing on the skateboard without introducing non-existent details.

Query: Is there a chair in the image?LLaVA: Yes, there is a white chair in the image.LLaVA: Yes, there is a white chair in the image.LLaVA w/ MARINE: No, there is nochair in the image.The only object present is a white bird.Query: What is the position of the skateboard in the image?LLaVA: The skateboard is on the ground, with the rider
standing on it and holding onto the trucks.LLaVA w/ MARINE: The skateboard is on the ground,
with the person standing on it.

Figure 4: Examples of hallucination mitigation by our proposed MARINE across multiple tasks: POPE on the GQA dataset, LLaVA-QA90 task on the MSCOCO dataset, and image captioning. Hallucinated objects generated by LVLM are highlighted in red.

499

500

511

512

513

514 515

516

517

518

519

520 521

6 CONCLUSIONS, LIMITATIONS AND FUTURE WORK

526 In this paper, we introduced a training-free and API-free framework MARINE to mitigate object 527 hallucination in LVLMs during its text generation process. Leveraging a pre-trained object grounding 528 vision encoder for a novel guidance framework in the multi-modal setting, MARINE effectively 529 and cost-efficiently reduces the hallucinations of five widely-used LVLMs, as assessed by various 530 metrics across different tasks. The inherent compatibility of the MARINE with various vision models and projection functions further underscores its flexibility. In contrast to post-generation correction 531 methods, MARINE strikes a balance between efficiency, instruction-following ability and effectiveness 532 in reducing object hallucinations. 533

Limitations and future work. While MARINE has demonstrated impressive performance by utilizing guidance from image-grounded models, there remains potential for further improvement through the integration of advanced aggregation methods, such as multi-agent debate (Du et al., 2023), into the MARINE framework. Additionally, although MARINE is specifically designed to mitigate object hallucination, which is the most significant issue in LVLMs, extending its application to address other types of hallucinations in both LLMs and LVLMs across a broader range of benchmarks would be highly advantageous.

540 ETHICS STATEMENT 541

This paper introduces research aimed at advancing the field of Large Language Models. We are confident that our work will contribute to significant social benefits, particularly by enhancing the accountability of LLMs through the reduction of hallucinatory outputs. Our proposed method, MARINE, holds the potential to improve the fairness of LLM interactions by effectively reducing biased hallucinations. To the best of our knowledge, we have not identified any negative effects associated with our research that merit highlighting in this discussion.

548 REPRODUCIBILITY STATEMENT

We provide detailed descriptions of our experimental setups, datasets, models, codes in the supplementary materials to ensure the reproducibility of MARINE. The full experimental settings and hyperparameters are presented in Appendix C.

553 REFERENCES

554

566

567

568

569

577

578

579

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022. 1
- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propositional
 image caption evaluation, 2016. 8, 16, 21
- Yakoub Bazi, Mohamad Mahmoud Al Rahhal, Laila Bashmal, and Mansour Zuair. Vision–language model for visual question answering in medical imagery. *Bioengineering*, 10(3):380, 2023.
- 564 Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing* 565 *text with the natural language toolkit.* "O'Reilly Media, Inc.", 2009. 4
 - Ali Furkan Biten, Lluís Gómez, and Dimosthenis Karatzas. Let there be a clock on the beach: Reducing object hallucination in image captioning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1381–1390, 2022. **3**, 4
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey
 Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pp. 213–229. Springer, 2020. 2, 5
- Fredrik Carlsson, Joey Öhman, Fangyu Liu, Severine Verlinden, Joakim Nivre, and Magnus Sahlgren.
 Fine-grained controllable text generation using non-residual prompting. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.
 6837–6857, 2022. 3
 - Pierre Chambon, Christian Bluethgen, Curtis P Langlotz, and Akshay Chaudhari. Adapting pretrained vision-language foundational models to medical imaging domains. *arXiv preprint arXiv:2210.04133*, 2022. 1
- Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023. 5, 16
- Zhaorun Chen, Zhuokai Zhao, Hongyin Luo, Huaxiu Yao, Bo Li, and Jiawei Zhou. Halc: Object hallucination reduction via adaptive focal-contrast decoding. *arXiv preprint arXiv:2403.00425*, 2024. 3, 16
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. 16
- 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, 2023a. 1, 3, 16

594 Wenliang Dai, Zihan Liu, Ziwei Ji, Dan Su, and Pascale Fung. Plausible may not be faithful: Probing 595 object hallucination in vision-language pre-training. In Proceedings of the 17th Conference of the 596 European Chapter of the Association for Computational Linguistics, pp. 2128–2140, 2023b. 3, 6 597 Yihe Deng, Pan Lu, Fan Yin, Ziniu Hu, Sheng Shen, James Zou, Kai-Wei Chang, and Wei Wang. 598 Enhancing large vision language models with self-training on image comprehension. arXiv preprint arXiv:2405.19716, 2024. 1, 22 600 601 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. Advances 602 in neural information processing systems, 34:8780–8794, 2021. 3, 4 603 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 604 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An 605 image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint 606 arXiv:2010.11929, 2020. 16 607 608 Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. arXiv preprint arXiv:2305.14325, 609 2023. 10 610 611 Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong 612 Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. 613 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 614 19358–19369, 2023. 16 615 Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, 616 Jinrui Yang, Xiawu Zheng, et al. Mme: A comprehensive evaluation benchmark for multimodal 617 large language models. arXiv preprint arXiv:2306.13394, 2023. 1, 3 618 619 Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, 620 Conghui He, Xiangyu Yue, Hongsheng Li, and Yu Qiao. Llama-adapter v2: Parameter-efficient 621 visual instruction model, 2023. 1, 3 622 Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and preventing hallucinations in large vision 623 language models. arXiv preprint arXiv:2308.06394, 2023. 1, 3 624 625 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In NeurIPS 2021 Workshop on 626 Deep Generative Models and Downstream Applications, 2021. 3, 4 627 Zhiting Hu and Li Erran Li. A causal lens for controllable text generation. Advances in Neural 628 Information Processing Systems, 34:24941–24955, 2021. 3 629 630 Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming 631 Zhang, and Nenghai Yu. Opera: Alleviating hallucination in multi-modal large language models 632 via over-trust penalty and retrospection-allocation. arXiv preprint arXiv:2311.17911, 2023a. 3, 6, 633 9.16.17 634 Xinyu Huang, Yi-Jie Huang, Youcai Zhang, Weiwei Tian, Rui Feng, Yuejie Zhang, Yanchun Xie, 635 Yaqian Li, and Lei Zhang. Open-set image tagging with multi-grained text supervision, 2023b. 636 URL https://arxiv.org/abs/2310.15200.2,5 637 638 Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning 639 and compositional question answering. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 6700-6709, 2019. 2, 6, 8, 21 640 641 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, 642 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. ACM 643 Computing Surveys, 55(12):1-38, 2023. 1 644 645 Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with 646 noisy text supervision. In International Conference on Machine Learning, pp. 4904–4916. PMLR, 647 2021. 1

658

659

660

- Liqiang Jing, Ruosen Li, Yunmo Chen, Mengzhao Jia, and Xinya Du. Faithscore: Evaluating hallucinations in large vision-language models. *arXiv preprint arXiv:2311.01477*, 2023. 1
- 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
 decoding. *arXiv preprint arXiv:2311.16922*, 2023. 3, 6, 9, 16, 17
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023a. 16
 - Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-Im improves controllable text generation. *Advances in Neural Information Processing Systems*, 35: 4328–4343, 2022. 3
- Kiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In
 Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers),
 pp. 4582–4597, 2021. 3
- 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*, 2023b. 1, 2, 3, 6, 18, 21
- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization* Branches Out, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL https://aclanthology.org/W04–1013. 8, 16, 21
- Sheng-Chieh Lin, Minghan Li, and Jimmy Lin. Aggretriever: A simple approach to aggregate
 textual representations for robust dense passage retrieval. *Transactions of the Association for Computational Linguistics*, 11:436–452, 2023a. 4, 5
- 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. 2, 6, 8, 21
- Zhaojiang Lin, Andrea Madotto, Yejin Bang, and Pascale Fung. The adapter-bot: All-in-one controllable conversational model. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 16081–16083, 2021. 3
- ⁶⁸³ Zhenghao Lin, Yeyun Gong, Yelong Shen, Tong Wu, Zhihao Fan, Chen Lin, Nan Duan, and Weizhu
 ⁶⁸⁴ Chen. Text generation with diffusion language models: A pre-training approach with continuous
 ⁶⁸⁵ paragraph denoise. In *International Conference on Machine Learning*, pp. 21051–21064. PMLR,
 ⁶⁸⁶ 2023b. 3
- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023a. 1
- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Mitigating
 hallucination in large multi-modal models via robust instruction tuning, 2023b. 1, 3
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023c. 5, 16
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023d. 1, 2, 3, 5, 6, 16, 19
- Shi Liu, Kecheng Zheng, and Wei Chen. Paying more attention to image: A training-free method for alleviating hallucination in lvlms. *arXiv preprint arXiv:2407.21771*, 2024. 3, 16
- Holy Lovenia, Wenliang Dai, Samuel Cahyawijaya, Ziwei Ji, and Pascale Fung. Negative object presence evaluation (nope) to measure object hallucination in vision-language models. *arXiv* preprint arXiv:2310.05338, 2023. 1, 3

702 703 704	Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Neural baby talk. In 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7219–7228, 2018. doi: 10.1109/CVPR.2018.00754. 18
705 706 707 708	Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts, 2024. 1
709 710 711 712	Aman Madaan, Amrith Setlur, Tanmay Parekh, Barnabás Poczós, Graham Neubig, Yiming Yang, Ruslan Salakhutdinov, Alan W Black, and Shrimai Prabhumoye. Politeness transfer: A tag and generate approach. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational</i> <i>Linguistics</i> , pp. 1869–1881, 2020. 3
713 714	Tong Niu and Mohit Bansal. Polite dialogue generation without parallel data. <i>Transactions of the Association for Computational Linguistics</i> , 6:373–389, 2018. 3
715 716 717 718 719	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35: 27730–27744, 2022. 3
720 721 722	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pp. 311–318, 2002. 8, 16, 21
723 724 725 726	Shrimai Prabhumoye, Alan W Black, and Ruslan Salakhutdinov. Exploring controllable text genera- tion techniques. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pp. 1–14, 2020. 3
727 728 729 730	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 <i>International conference on machine learning</i> , pp. 8748–8763. PMLR, 2021. 1, 4, 16
731 732 733	Leonardo FR Ribeiro, Yue Zhang, and Iryna Gurevych. Structural adapters in pretrained language models for amr-to-text generation. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pp. 4269–4282, 2021. 3
734 735 736 727	Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pp. 4035–4045, 2018. 2, 3, 6, 18
738 739 740	Guillaume Sanchez, Honglu Fan, Alexander Spangher, Elad Levi, Pawan Sasanka Ammanamanchi, and Stella Biderman. Stay on topic with classifier-free guidance. <i>arXiv preprint arXiv:2306.17806</i> , 2023. 3, 4, 5
741 742 743	Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A benchmark for visual question answering using world knowledge. In <i>European</i> <i>Conference on Computer Vision</i> , pp. 146–162. Springer, 2022. 2, 6, 8, 21
744 745 746 747	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023. 16
748 749	Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation, 2015. 8, 16, 21
750 751 752 753	David Wan, Jaemin Cho, Elias Stengel-Eskin, and Mohit Bansal. Contrastive region guidance: Improving grounding in vision-language models without training. <i>arXiv preprint arXiv:2403.02325</i> , 2024. 3 , 16
754 755	Bin Wang, Fan Wu, Xiao Han, Jiahui Peng, Huaping Zhong, Pan Zhang, Xiaoyi Dong, Weijia Li, Wei Li, Jiaqi Wang, et al. Vigc: Visual instruction generation and correction. <i>arXiv preprint arXiv:2308.12714</i> , 2023a. 1, 3

756 757 758 750	Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. Git: A generative image-to-text transformer for vision and language. <i>arXiv</i> preprint arXiv:2205.14100, 2022. 23
760 761 762	Junyang Wang, Yiyang Zhou, Guohai Xu, Pengcheng Shi, Chenlin Zhao, Haiyang Xu, Qinghao Ye, Ming Yan, Ji Zhang, Jihua Zhu, et al. Evaluation and analysis of hallucination in large vision-language models. <i>arXiv preprint arXiv:2308.15126</i> , 2023b. 1, 3
763 764 765	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In <i>The Eleventh International Conference on Learning Representations</i> . 5
766 767 768 769	Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models. <i>arXiv preprint arXiv:2306.09265</i> , 2023. 1
770 771 772 773	Zhengyuan Yang, Jianfeng Wang, Zhe Gan, Linjie Li, Kevin Lin, Chenfei Wu, Nan Duan, Zicheng Liu, Ce Liu, Michael Zeng, et al. Reco: Region-controlled text-to-image generation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 14246–14255, 2023. 3
774 775 776	Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. <i>arXiv preprint arXiv:2304.14178</i> , 2023. 1, 3, 5, 16
778 779 780	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. 1, 3, 6, 7, 9, 18, 28
781 782 783 784	Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun, et al. Rlhf-v: Towards trustworthy mllms via behavior alignment from fine-grained correctional human feedback. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 13807–13816, 2024. 1
785 786 787 788	Bohan Zhai, Shijia Yang, Chenfeng Xu, Sheng Shen, Kurt Keutzer, and Manling Li. Halle-switch: Controlling object hallucination in large vision language models. <i>arXiv e-prints</i> , pp. arXiv–2310, 2023. 3, 9
789 790 791	Hanqing Zhang, Haolin Song, Shaoyu Li, Ming Zhou, and Dawei Song. A survey of controllable text generation using transformer-based pre-trained language models. <i>ACM Computing Surveys</i> , 56(3): 1–37, 2023a. 3
792 793 794	Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Peng Gao, et al. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? <i>arXiv preprint arXiv:2403.14624</i> , 2024a. 1
795 796 797 798	Yuechen Zhang, Shengju Qian, Bohao Peng, Shu Liu, and Jiaya Jia. Prompt highlighter: Interactive control for multi-modal llms. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 13215–13224, 2024b. 3, 16
799 800 801	Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal chain-of-thought reasoning in language models, 2023b. 2, 4
802 803 804	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. 1, 3, 6, 9, 16, 17, 28
805 806 807	Yiyang Zhou, Chenhang Cui, Rafael Rafailov, Chelsea Finn, and Huaxiu Yao. Aligning modalities in vision large language models via preference fine-tuning. <i>arXiv preprint arXiv:2402.11411</i> , 2024. 1
808 809	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. 1 , 3

810 A BROADER IMPACT STATEMENT

By mitigating hallucinations, MARINE has the potential to offer a positive social impact by ensuring that LVLMs generate more accountable responses. Despite this merit, MARINE cannot address prejudicial biases inherent in LLM prior knowledge, which could be a focus of future work. Aside from this limitation, we do not foresee any potential negative impacts of our work that need to be specifically highlighted at this time, to the best of our knowledge.

817 818

819

836

837

838 839

846

847

848

849

850

851

852

853

854

855

856

858

B CONCURRENT WORKS.

Concurrently, Leng et al. (2023) introduced Visual Contrastive Decoding (VCD), a technique that 820 applies noise to image inputs and penalizes logit outputs of these corrupted images. Huang et al. 821 (2023a) enhanced beam-search decoding with the Over-trust Penalty and Retrospection-Allocation 822 Strategy (OPERA), which penalizes over-trust and refines token selection based on previous outputs. 823 HALC (Chen et al., 2024) employs adaptive focal-contrast decoding to encourage LVLMs to focus 824 on fine-grained visual information, while using a computationally intensive beam search algorithm. 825 Liu et al. (2024) addresses the issue of "Text Inertia", where LVLMs produce the same hallucinated 826 object descriptions regardless of whether an image is actually provided. Wan et al. (2024) focuses on 827 fine-grained guidance and especially operates at the sub-image level. Zhang et al. (2024b) adjusts 828 the attention weights of the visual tokens within the attention mechanism to control the focus of the 829 generation on the highlighted parts.

830 C EXPERIMENT DETAILS

We conduct all of the experiments using 8 A6000 GPU with 48GB GPU memory. Each single experiment can be run on a single A6000 GPU.

834 C.1 MODEL ARCHITECTURES

In Table 8, we provide detailed descriptions of the LVLM architectures used in our experiments. These LVLMs respectively leverage the pre-trained vision encoder of the models we listed, which are all based on the Vision Transformer (ViT) (Dosovitskiy et al., 2020) architecture.

fuble 0. Details of the Dy Ent areniteetares that we ased in our paper.

340	Model	Vision encoder	LLM
841	LLaVA (Liu et al., 2023d)	CLIP-L (Radford et al., 2021)	LLaMA-2-7B-Chat (Touvron et al., 2023)
842	LLaVA-v1.5 (Liu et al., 2023c)	CLIP-L-336px (Radford et al., 2021)	Vicuna-v1.5-7B (Chiang et al., 2023)
8/13	MiniGPT-v2 (Chen et al., 2023)	EVA-G (Fang et al., 2023)	LLaMA-2-7B-Chat (Touvron et al., 2023)
040	mPLUG-OWL2 (Ye et al., 2023)	CLIP-L (Radford et al., 2021)	LLaMA-2-7B (Touvron et al., 2023)
844	InstructBLIP (Dai et al., 2023a)	BLIP-2 (Li et al., 2023a)	Vicuna-v1.1-7B (Chiang et al., 2023)
845			

C.2 DESCRIPTIONS ABOUT ADDITIONAL METRICS

In Figure 2, we evaluate the text quality of the outputs generated with MARINE using general metrics as follows:

- *BLEU* (Papineni et al., 2002) measures how well the generated translation matches the reference translations in terms of n-gram overlap.
- *ROUGH-L* (Lin, 2004) measures the quality of a machine-generated summary by comparing it to one or more reference summaries.
- *CIDEr* (Vedantam et al., 2015) assesses the quality of image captioning models. It focuses on evaluating how well the generated captions align with human consensus.
- *SPICE* (Anderson et al., 2016) focuses on assessing the semantic similarity between the generated captions and reference captions.
- C.3 PROMPT TEMPLATES

For each query, we randomly select a prompt template from the available template list, as shown in Table 9.

- 861 C.4 DETAILS OF BASELINES
- Specifically, the hyperparameters for LURE (Zhou et al., 2023), VCD (Leng et al., 2023), OPERA (Huang et al., 2023a) are reported in Table 10, 11 and 12 respectively. We strictly

864	Table	9: Details of the LVLM architectures that we used in our paper.
865	Template Type	Prompt Template
866 867	MARINE-intersec	This image contains <object_grounding>. Based on this, <query> The image contains the following objects: <object_grounding>. Given these datacted objects: <ouery></ouery></object_grounding></query></object_grounding>
000		This image shows the following objects: <object grounding="">. Using this infor-</object>
009		mation, <query></query>
070		The objects found in this image are: <object_grounding>. Considering this list</object_grounding>
071	DODE took	of objects, <query></query>
072	r Or E task	assume anything beyond these objects. Based solely on this list. <0UERY>
073		The detected objects in the image are: <object_grounding>. Answer based only</object_grounding>
874		on these objects. <query></query>
875		This image shows the following objects: <object_grounding>. You must answer</object_grounding>
876		Using only the objects in this list. Given these detected objects, <query></query>
877		rely strictly on this list of objects and make no other guesses. Based on this, <ouery></ouery>
878	MARINE-union	List of detected objects in the image:
879		<object_grounding_a></object_grounding_a>
880		<object_grounding_b></object_grounding_b>
881		Based on the detected objects above, <query></query>
882		<pre></pre>
883		<pre><0BJECT_GROUNDING_B></pre>
884		Given these findings, <query></query>
885		The following objects were detected in the image:
886		<pre><object_grounding_a> </object_grounding_a></pre>
887		With this information <ouery></ouery>
888		Here is a list of all objects detected in the image:
889		<object_grounding_a></object_grounding_a>
890		<object_grounding_b></object_grounding_b>
891		Do not infer or hallucinate any additional objects. Using only the detected objects,
892		<pre></pre>
893	followed the original	l implementations and default hyperparameters described in their papers to
894	reproduce the results	for each baseline.
805	Ta	ble 10: LURE (Zhou et al., 2023) Hyperparameter Settings
806		Parameters Value
207		Uncertainty Threshold γ 0.9
808		Position Threshold ι 0.8
899	Т	able 11: VCD (Leng et al., 2023) Hyperparameter Settings
900		Parameters Value
901		Amplification Factor α 1
902		Adaptive Plausibility Threshold 0.1
903		Diffusion Noise Step 500
904	Tabl	e 12: OPERA (Huang et al., 2023a) Hyperparameter Settings
905		Parameters Value
906		Self-attention Weights Scale Factor θ 50
907		Attending Retrospection Threshold 25
908		Beam Size 5
909		Attention Candidates 1
910		Penalty Weights 1
911		
912		

C.5 EXPERIMENT SETTING FOR HALLUCINATION EVALUATIONS

815 Key factors that potentially affect the hallucination evaluation outcomes, including the evaluation
916 dataset and prompt template, LVLM's sampling strategy and batched generation techniques, and
917 guidance strength, are detailed in this section. The hyper-parameters setting for MARINE and overall experiment settings are shown in Table 13 and 14.

931

933 934 935

920	Parameters	Value
921	Guidance	, unue
922	Guidance Strength	0.7
923	Noise Intensity for DETR	0.95
924	Detect Threshold for RAM++	0.68
925	Generation	
926	Max Token Length	64
927	Sampling	Greedy
928	Random Seed	242

918 Table 13: MARINE Hyperparameter Settings. The settings are fixed depending on the question-919 answering tasks. 92

929 Table 14: Batch size for LVLM generation is fixed across all experiments unless otherwise noted. To 930 expedite the evaluation process, we employed the batched generation. We avoid the negative impact of batched generation by adopting left padding if the LVLM does not explicitly assign the padding strategy for inference. 932

Model	LLaVA	LLaVA-v1.5	MiniGPTv	mPLUG-Owl2	InstructBLIP
Batch Size	16	4	32	16	16

936 **Experiment setting for CHAIR evaluation.** We adopt the same prompt "Generate a short caption 937 of the image." as utilized by Li et al. (2023b). The hyperparameters are fixed, including a guidance 938 strength of 0.7, noise intensity for DETR at 0.95, a detection threshold for RAM++ of 0.68, a 939 maximum token length of 64, and a greedy sampling approach with a random seed of 242.

940 For the calculation of CHAIR metrics, we referenced the 80 object categories annotated in the 941 MSCOCO dataset, following Rohrbach et al. (2018). Besides, we employed the synonym list from 942 Lu et al. (2018) to align synonymous words in the generated text with MSCOCO object categories. Additionally, due to the cost considerations associated with the GPT-3.5 API, we limited our analysis 943 to 200 samples for Woodpecker correction for each model and reported the result in Table 1. 944

945 **Experiment setting for POPE evaluation.** POPE is a flexible approach to evaluating hallucinations 946 in LVLMs, which formulates a binary classification task by prompting LVLMs with questions such as 947 "Is there a keyboard in this image?" to answer "yes" or "no". Following Li et al. (2023b), we created 948 3000 POPE questions across three datasets-500 images each from MSCOCO, A-OKVQA, and 949 GQA for the POPE evaluation. We reported the adversarial settings in Table 2, the most challenging setting, which constructs POPE questions from the top-k most frequently co-occurring but absent 950 objects. Additionally, in Table 4, we reported the average scores under random, popular, adversarial 951 settings across MSCOCO, A-OKVQA, and GQA datasets. The full POPE results are in Tabel 15. 952

Similarly, we constrained our analysis to 200 samples for Woodpecker correction for each model due 953 to the high costs associated with the GPT API. The outcomes of this analysis are detailed in Table 2. 954

955 **Experiment setting for GPT-4V-aided evaluation.** The GPT-4V-aided evaluation compares the 956 outputs of two LVLM assistants using GPT-4V as a judge. We prompted GPT-4V to assess the quality 957 of the generated outputs, scoring them out of 10 in two aspects:

- 958 • Accuracy: how accurately each assistant describes the image;
- 959 • Detailedness: the richness of necessary details in the response.

960 As shown in Figure 5, the assessment prompt template we used is slightly different from that 961 of Yin et al. (2023). Specifically, we include the original question for a task-orientated evalu-962 ation and exclude prompts that describe Woodpecker-specific output formats like object bounding boxes. Examples of the GPT-4V-aid evaluation responses are illustrated in Figure 6 and 7. 963 Besides, a fixed guidance strength of 0.5 was used in the evaluations in Table 3. Utilizing the 964 gpt-4-1106-vision-preview, all final experiments were conducted between 01/01/2024-965 01/30/2024. As OpenAI continues to update its API, accessing different versions may result in 966 slightly different values. 967

968 **Experiment setting for ablation study.** To explore different methods of integrating imagegrounding models, we investigate the intersection and union of detected objects, with integration 969 based on synonyms using the NLTK package. 970

To quantitatively assess the influence of guidance strength, we varied it from 0 to 1, as shown in 971 Figure 10. These quantitative experiments were conducted using the same setting as those in CHAIR

evaluation. For qualitative analysis, exemplified in Figure 13 and 10, we selected guidance strength from a recommended range of $\gamma \in (0.3, 0.7)$.

Prompt 975 You are required to score the performance of two AI assistants in describing a given image. You should pay extra attention to the 976 hallucination, which refers to the part of descriptions that are inconsistent with the image content, such as claiming the existence of something not present in the image 977 Please rate the responses of the assistants on a scale of 1 to 10, where a higher score indicates better performance, according to the 978 following criteria 1: Accuracy: whether the response is accurate with respect to the image content. Responses with fewer hallucinations should be 979 given higher scores 2: Detailedness: whether the response is rich in necessary details. Note that hallucinated descriptions should not count as necessary 980 details 981 Please output a single line for each criterion, containing only two values indicating the scores for Assistant 1 and 2, respectively. The two scores are separated by a space. Following the scores, please provide an explanation of your evaluation, avoiding any 982 potential bias and ensuring that the order in which the responses were presented does not affect your judgment.". 983 Please score the performance of two AI assistants in describing a given image following the given question 984 Question: 985 {question} 986 Assistant 1: 987 {answer 1} 988 Assistant 2: 989 {answer 2} 990 Output format: Accuracy 991 Scores of the two answers: 992 Reason: 993 Detailedness: Scores of the two answers: 994 Reason: 995

Figure 5: Prompt template for GPT-4V-aided evaluation. {question} is the original instruction;
 {answer 1} is the original response, and {answer 2} is the response generated by the LVLM using
 MARINE-DETR with a guidance strength of 0.5.

999 1000 C.6 Experiment Setting on Other Vision-Language Tasks

Experiment setting for text quality analysis. For text quality analysis, we adopted 90 visual questions from the LLaVA-QA90 task (including conversations, visual perceptions, and complex reasoning subtasks), and randomly selected 50 MSCOCO images for image captioning task. Following Liu et al. (2023d), we adpoted the response generated by text-only GPT-4 (0314) with the context captions/boxes provided. answers given by GPT-4 as references for LLaVA-QA90 task and used image captions provided in MSCOCO annotations as references for image captioning task.

In Table 16 and Table 17, we present a detailed evaluation on the image captioning task for both
 MSCOCO and LLaVA-QA90 using metrics including BLEU, ROUGE, CIDEr and SPICE. The
 corresponding figure result is shown in Figure 2.

Experiment setting for latency analysis. We compared our method with existing baselines in terms of the trade-off between inference cost and the effectiveness of reducing object hallucinations, as shown in Table 5. For post-correction baselines such as Woodpecker and LURE, we first prompted LLaVA (llava-llama-2-7b-chat-lightning-preview) to generate captions and then measure the latency of generating the corrected outputs. The total latency for post-correction baselines includes both the generation and correction processes. For decoding methods such as VCD, OPERA and our method, we measured the latency of LLaVA generating captions directly.

We prompted the models with "Generate a short caption of the image." on 500 MSCOCO images with a batch size of 1 and a maximum token length of 64, without any stopping criteria, using a single A6000 GPU. Then latency was calculated as the ratio of the number of output tokens and encoding and generation time.

1021

1000

1022

1023

1024

Prompt
Question: Can you eleborate on the elements of the picture provided?
Assistant 1 (LLaVA-v1.5): The image features a desk with various items on it. There is a bookshelf with a few books.
including a prominent one titled \"Administrative Law \" A blue and white tie is placed on the desk,
adding a touch of color and formality to the scene. A pair of sunglasses is also resting on the desk, possibly belonging to someone who was recently there.
In addition to the books and the there is a sum and an analy on the dark indicating that the summa
might have been enjoying a snack or drink while working or studying. A teddy bear is also present
in the scene sitting on the desk, and adding a playful touch to the setting.
Assistant 2 (LLaVA-v1.5 w/ MARINE):
The image features a desk with various items on it. There is a book, possibly a law book, positioned on the left side of the desk. A blue and white tie is placed next to the book, adding a
touch of color to the scene. On the right side of the desk, there is an apple, possibly a Granny
Smith apple, and a pair of sungrasses.
Additionally, there is a stuffed teddy bear located at the top right corner of the desk, giving the
owner might have an interest in sports.
GPT-4V Evaluation
Scores of the two answers: 4 7
Reason: Both assistants mention items that are present in the image, such as the book with the title "Administrative Law," the tie, the apple the supplayers and the teddy hear. Assistant I hallucinates the presence of a cup, which is not visible in the image
Assistant 2 inaccurately describes the presence of a baseball glove, which is not in the image; the item mistaken for a glove is
actually part of the teddy bear's outfit. Assistant 2 also avoids mentioning the specific title of the book, even though it is visible in the image.
Detailedness: Scores of the two answers: 5 6
Reason: Assistant 1 provides a more detailed description of the scene by elaborating on the potential reasons for the items' presence,
such as the person enjoying a snack of the teddy bear adding a playful fouch, However, the mention of a cup, which is not in the image, detracts from the quality of the response. Assistant 2 provides a sufficient amount of detail without over-elaborating and
makes a subject of the second state of the sec
makes a specific guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant I result in a slightly lower score.
makes a specific guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score.
Figure 6: An example of GPT-4V-aided evaluation on a detailed question.
Figure 6: An example of GPT-4V-aided evaluation on a detailed question.
rest a spectric guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question.
Prompt Question: What might be the purpose of this table arrangement?
Prompt Question: What might be the purpose of this table arrangement?
reaces a spectre guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prompt Question: What might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace.
Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Frompt Question: What might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading.
Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Frompt Question: What might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation and snacking. The presone of the book and sunglasses suggests that the person might be taking a break from reading or working on a project, while the apple provides a healthy snack
Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Frompt Question: What might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation and snacking. The presence of the book and sunglasses suggests that the person might be taking a break from reading or working on a project, while the apple provides a healthy snack option. The combination of these items on the table indicates that the person values both work and an apple area to the table indicates that the person values both work and are relaxation and snacking.
makes a spectre guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prompt Question: What might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation and snacking. The presence of the book and sunglasses suggests that the person might be taking a break from reading or working on a project, while the apple provides a healthy snack option. The combination of these items on the table indicates that the person values both work and relaxation, aiming to create a balance between the two.
makes a spectre guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prompt Question: What might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation and snacking. The presence of the book and sunglasses suggests that the person might be taking a break from reading or working on a project, while the apple provides a healthy snack option. The combination of these items on the table indicates that the person values both work and relaxation, aiming to create a balance between the two. Assistant 2 (LLaVA-v1.5 w/MARINE): The purpose of this table arrangement appears to be a combination of work and leisure. The
makes a spectre guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prompt Question: What might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation and snacking. The presence of the book and sunglasses suggests that the person might be taking a break from reading or working on a project, while the apple provides a healthy snack option. The combination of these items on the table indicates that the person values both work and relaxation, aiming to create a balance between the two. Assistant 2 (LLaVA-v1.5 w/ MARINE): The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and
makes a spectre guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prompt Question: What might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation and snacking. The presence of the book and sunglasses suggests that the person might be taking a break from reading or working on a project, while the apple provides a healthy snack option. The combination of these items on the table indicates that the person values both work and relaxation, aiming to create a balance between the two. Assistant 2 (LLaVA-v1.5 w/ MARINE): The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears ato be a combination of work an
makes a spectre guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prompt Question: What might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation and snacking. The presence of the book and sunglasses suggests that the person might be taking a break from reading or working on a project, while the apple provides a healthy snack are laxation, aiming to create a balance between the two. Assistant 2 (LLaVA-v1.5 w/ MARINE): The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appeares appeare and work are study. Additionally, the a
<section-header>makes a spectre guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prompt Question: Water of the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, potential or of the sole and sunglasses suggests that the person might be taking a break from reading or working on a project, while the pape provides a healthy snace are lavation, aiming to create a balance between the two. Assistant 2 (LLaVA-v1.5 w/ MARINE): The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement</section-header>
<section-header> makes a spectre guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prompl Question: Mat might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation and snacking. The presence of the book and sunglasses suggests that the person might be table arrangement apples to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of the table and an apple suggests that the person values both work and the person study, as indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indi</section-header>
<section-header>makes a spectre guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prompl Question: Wat might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, pair abreak from reading or working on a project, while the apple provides a heat person might be table indicates that the person values both work and relaxation, aiming to create a balance between the two. Assistant 2 (LLaVA-v1.5 w/ MARINE): The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of this table arrangement appears to be a combination of work and leisure. The purpose of the table taking a break, enjoying a snack, or preparing for a more relaxed at the person might be taking a break, enjoying a snack, or preparing for a more relaxed at the person appear of the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally, the apple and sunglasses are indicated by the books. Additionally the apple and sunglasses are indicated by the books. Additionally the apple and sunglasses are indicated by the books. Additionally the apple and sunglasses are indicated by the books. Additionally the apple and sunglasses are indicated by the books. Additionally the apple and sunglasses are indicated by the books. Additionally the ap</section-header>
makes a spectre guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prompl Question: What might be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation and snacking. The presence of the book and sunglasses suggests that the person wing the sequence of the stable arrangement apple in details about the type of this table arrangement apple. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation and snacking. The presence of the book and sunglasses suggests that the person values both work and relaxation, aiming to create a balance between the two. Assistant 2 (LLaVA-v1.5 w/ MARINE): The purpose of this table arrangement appears to be a combination of work and leisure. The presence of the teddy bear adds a touch of playfulness and personalization to the setting, creating a cozy and inviting atmosphere. CPT-4V Evaluation Accuracy: Cores of the two answers: 67 Resons the wo answers: 67 Resons the wo answers: 67 Resons of the sub and a represention, she sitat 1 incorrectly states that there's a book on the table, which
makes a spectric guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prompt Question: What might be the purpose of this table arrangement? Assistant 1 (LaVA-v1.5): The purpose of this table arrangement might be to create a comfortable and productive workspace. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading. The table features a pair of sunglasses, a book, and an apple, which are essential items for reading. The purpose of this table arrangement might be to create a comfortable and productive workspace. The sentence of the table areangement meet the book and sunglasses suggests that the person might be taking a break, reinoging a snack, or preparing for a more relaxed activity. The presence of the teddy bear adds a touch of playfulness and personalization to the sentence areands and the taking a break, enjoying a snack, or preparing for a more relaxed activity. The presence of the teddy bear adds a touch of playfulness and personalization to the sentence meets and inviting atmosphere. CPT-4V Evaluation Accuracy: Scores of the two answers: 6.7 Reason: Both Assistant and a few errors in their descriptions, Assistant 1 incorrectly states that the ters's a book on the table, which is atmosphere. Scores of the two answers: 6.7 Reason: Both Assistant and a few errors in their descriptions. Assistant 1 incorrectly states that the ters's
<text><section-header><text><text><text><text><text><text><text></text></text></text></text></text></text></text></section-header></text>
<section-header><section-header>Bases a spectra guess about me type or apple. However, the incorrect mention of a baseball glove and the lack of details about the look titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Form 1 Westion: What might be the purpose of this table arrangement? Assistant 1 (LaVA-v1.5) The table features a pair of sunglasses, a book, and an apple, which are essential items for reading the table areangement relaxion and snacking. The presence of the book and sunglasses suggests that the person might be table assistant 5 a society of the society of the society of the society of the suble arrangement apple is look and sunglasses suggests that the person might be table assistent to the person values both work and a preak form reading or working on a project, while the apple provides a healthy snach apple and sunglasses at the table indicates that the person values both work and a lange to the society of the suble arrangement appears to be a combination of work and leisure. The presence of the table by the books. Additionally, the apple and sunglasses at the terison of book, a sunglasses, at the outch of playfulness and personalization ot be stiring, creating a cozy and inviting atmosphere. CPT-4V Evaluation Accuracy: CPT-4V Evaluation Exoson for the table area form forms in the indescriptions, Assistant 1 incorrectly states that the relax to in and the tie is in mentioned, which is a significant element in the image. Assistant 2 correctly recognizes the presence of books in a sunglasses, at the subs, area not on the table but more likely on a shelf next to it and the tie is in mentioned, which is a significant element in the image. Assistant 2 correctly recognizes the presence of books, a sunglasses, at the sunglasses, at its sunglasses, and an apple,</section-header></section-header>
<text><section-header><section-header><section-header><section-header><section-header><section-header><section-header></section-header></section-header></section-header></section-header></section-header></section-header></section-header></text>
<section-header>Because a spectre guess arout the type of apple, However, the incorrect mention of a baseball glove and the lack of details about the lock titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Prove the model of the stable arrangement of the stable arrangement? And might be the purpose of this table arrangement might be to create a comfortable and productive workspace fraction and snacking. The presence of the book and sunglasses suggests that the person neight of sunglasses, a book, and an apple, which are essential items for neight of the table features a pair of sunglasses, a book, and an apple, which are essential items for neight of the table features a pair of sunglasses, a book, and an apple, which are essential items for neight of the table features a pair of sunglasses, a book, and an apple, which are essential items for neight of the stable arrangement might be to create a comfortable and productive workspace relaxation, aining to create a balance between the two. Assiant 2 (LaVA-v1.5 w/ MARINE) The propose of this table arrangement appears to be a combination of work and leisure. The be engaged in work or study, as indicated by the books. Additionally, the apple and sunglasses are sunglasses. A site, and an apple suggests that the person using this space may be engaged in work or study, as indicated by the books. Additionally, the apple and sunglasses are surged a court of playfulness and personalization to the study. Streame and the study a brooks. Additionally, the apple and sunglasses are surger and the table indicates that the person using this space may be engaged in work or study, as indicated by the books. Additionally, the apple and sunglasses. A book on the table, which are streamed are create a core of the tody has a touch of playfulness and personalization to the streame and in viting atmosphere. CPT-4V Evaluation Secone of the two answers: 67 Secon</section-header>
Index a spectric guess about the type of apple. However, the incorrect mention of a baseball glove and the lack of details about the book titles compared to Assistant 1 result in a slightly lower score. Figure 6: An example of GPT-4V-aided evaluation on a detailed question. Frompt Question: What night be the purpose of this table arrangement? Assistant 1 (LLaVA-v1.5) The purpose of this table arrangement might be to create a comfortable and productive workspace. The purpose of this table arrangement might be to create a comfortable and productive workspace. The purpose of this table arrangement might be to create a comfortable and productive workspace. The about of the subale features a pair of sunglasses, a book, and an apple, which are essential items for reading, relaxation, aiming to create a balance between the two. Assistant 2 (LLaVA-v1.5 w/ MARINE): The prosence of this table arrangement tappears to be a combination of work and leisure. The be engaged in work or study, as indicated by the books. Additionally, the apple and sunglasses and and angles, a sunglasses, a tile, and an apple suggests that the person using this space may are its ation, aiming to create a balance between the two. Assistant 2 (LLaVA-v1.5 w/ MARINE): The prosence of books, a sunglasses, a tile, and an apple suggests that the person using this space may are its ation, a work, as indicated by the books. Additionally, the apple and sunglasses indicate that the person might be taking a break, enjoying a snack, or preparing for a more relaxed arity are creating a cozy and inviting atmosphere. CPT-VE Evaluation Meximal Scores Of the two answers: 67 Scores Boh kasistants made a few errors in their descriptions, Assistant 1 incorrectly states that there/s a book on the table, which is is asignificant element in the image. Assistant 2 correctly recognizes the presence of books, a tie, sunglasses, and an apple, along with is ambigous; technically, the books are not on the tabble turnore likely on a shelf next to i
<text><section-header><section-header><text><section-header><text><text><text><text></text></text></text></text></section-header></text></section-header></section-header></text>
<text><section-header><section-header><section-header><section-header><section-header><table-container><section-header></section-header></table-container></section-header></section-header></section-header></section-header></section-header></text>
<text><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header></section-header></section-header></section-header></section-header></section-header></section-header></section-header></section-header></section-header></text>

Dataset	Туре	Model	w/MARINE	Accuracy \uparrow	Precision \uparrow	Recall \uparrow	F1 ↑	Yes(%)
		LL oVA	X	51.8	50.9	99.5	67.4	97.7
	Adversarial	LLaVA	\checkmark	66.9	61.7	89.1	72.9	72.3
	/ tuversuriai	mPLUG-Owl2	×	72.5	65.5	94.9	77.5	72.4
		III 200 0 112	\checkmark	82.8	83.4	82.0	82.7	49.2
		LLaVA	×	52.4	51.2	99.8	6/./	97.4
MSCOCO	Popular		\checkmark	71.3	65.8	88.9	75.6	67.5
	•	mPLUG-Owl2	^	/5.8	08.7	94.9 82.0	/9./ 95.1	09.0 46.4
			×	58.3	00.4 54 5	02.0 00 7	85.1 70.5	40.4 01 /
	Random	LLaVA	Ś	78.5	73.4	89.3	80.6	60.8
			x	81.8	75.2	94.9	83.9	63.1
		mPLUG-Owl2	\checkmark	88.1	93.4	81.9	87.3	43.9
		LL aVA	X	50.0	50.0	99.5	66.6	99.5
	Advargial	LLaVA	\checkmark	56.3	53.6	94.3	68.3	88.1
	Auversiai	mPLUG-Owl2	×	62.5	57.3	98.1	72.3	85.6
	Popular		\checkmark	74.4	68.8	89.3	77.7	64.9
		LLaVA	×	50.1	50.1	99.8	66.7	99.7
A-OKVOA			\checkmark	63.0	58.0	94.5	71.9	81.6
·····		mPLUG-Owl2	×	69.1	62.1	97.9	76.0	78.9
			~	82.5	78.8	89.1	83.6	56.5
		LLaVA	×	55.4 72.7	52.8	99.8	69.1	94.4
	Random		v v	13.1	00.7 60.2	94.7	/8.3	71.0
		mPLUG-Owl2	Ŷ	80.2	80.2	98.2 80.3	80.2	71.0 50.1
			• •	59.2	50.1	09.5	69.2	30.1
		LLaVA	×	50.3	50.1	99.8	66.8	99.5
	Adversial		√ ▼	54.4	52.5	93.8	0/.3	89.4
		mPLUG-Owl2	<u>^</u>	76.0	73.6	90.2 81.2	75.0	19.0
			x	50.1	50.0	99.8	66.7	99.7
		LLaVA	, ,	58.7	55.1	94.3	69.5	85.5
GQA	Popular		×	70.6	63.8	94.9	76.3	74.4
		mPLUG-Owl2	\checkmark	77.6	75.6	81.3	78.4	53.8
		T T . 374	×	55.7	53.0	99.8	69.2	94.1
	Dondom	LLavA	\checkmark	74.3	67.3	94.8	78.7	70.5
	Kandom	mPLUG Owl2	×	82.0	75.2	95.5	84.1	63.5
		III LUG-OWI2	\checkmark	86.8	91.5	81.3	86.1	44.4

1080	
	Table 15: Detailed POPE (L1 et al., 2023b) results on three datasets (MSCOCO (L1n et al., 2014),
1081	A-OKVOA (Schwenk et al., 2022), GOA (Hudson & Manning, 2019)).
1000	

1111Table 16: Performance on general metrics for the image captioning task, including BLEU (Papineni1112et al., 2002), ROUGE-L (Lin, 2004), CIDEr (Vedantam et al., 2015) and SPICE (Anderson et al.,11132016) scores(%).

Model	w/MARINE	BLEU_1 (†)	BLEU_2 (†)	BLEU_3 (†)	BLEU_4 (\uparrow)	ROUGE_L (†)	CIDEr (†)	SPICE (†
LLaVA	×	14.06	7.12	3.72	1.90	22.06	0.08	16.77
	√	18.59	9.96	5.47	3.04	26.02	0.21	20.58
mPLUG-Owl2	×	39.91	25.16	16.57	11.24	36.26	1.05	26.82
	√	39.51	24.37	15.93	10.70	36.01	1.03	27.42

Table 17: Performance on general metrics for the LLaVA-QA90 task, including BLEU (Papineni et al., 2002), ROUGE-L (Lin, 2004), CIDEr (Vedantam et al., 2015) and SPICE (Anderson et al., 2016) scores(%).

1122	Model	w/MARINE	BLEU_1 (†)	BLEU_2 (†)	BLEU_3 (†)	BLEU_4 (†)	ROUGE_L (†)	CIDEr (†)	SPICE (\uparrow)
1123	LLaVA	× √	21.02 23.37	12.91 14.39	8.79 9.59	6.41 6.83	32.30 33.81	0.93 0.99	31.36 31.91
124 125	mPLUG-Owl2	× √	44.50 45.82	28.57 28.87	19.58 19.24	14.43 13.70	40.24 38.54	1.46 1.29	40.51 38.70

D ADDITIONAL EXPERIMENTS

1128 D.1 ADDITIONAL BASELINES

1129 D.1.1 DIRECT PROMPTING 1130

1119

1120

1121

1126

1127

We conducted additional experiments to compare our approach with a baseline that uses carefullyengineered prompts designed to reduce hallucination:

1133 Describe the visible contents of this image in as much detail as possible without adding any information not clearly visible. Only mention objects, colors, shapes, and textures that can be directly

Method		LLaVA		1	LLaVA-v1.	5	n	nPLUG-Owl	2
CHAIR	$C_s\downarrow$	$C_i\downarrow$	Recall ↑	$C_s\downarrow$	$C_i\downarrow$	Recall ↑	$C_s \downarrow$	$C_i\downarrow$	Recall ↑
Original	26.6	10.5	47.4	8.8	4.6	41.1	5.0	3.2	33.2
Direct Prompting	27.2	11.0	46.4	19.6	8.3	52.3	9.0	5.1	42.0
Prompts as Additional Guidance	37.4	10.5	50.4	12.6	5.9	44.6	6.6	3.9	40.4
MARINE (ours)	17.8	7.2	50.8	6.2	3.0	44.3	4.2	2.3	41.4

Table 18: Comparison against carefully engineered prompts.

Table 19: Comparison against ensembling multiple LVLMs. We consider ensembling all possible combinations of the five LVLMs and report the average performance with the standard deviation.

Method	Accuracy \uparrow	F1 \uparrow	Yes Ratio
Voting on 2 LVLMs MARINE(ours)	$\begin{array}{c} 76.7\pm7.8\\ 79.9\pm7.4\end{array}$	$\begin{array}{c} 79.8\pm4.8\\ 80.4\pm4.6\end{array}$	$\begin{array}{c} 62.6 \pm 11.2 \\ 51.1 \pm 12.6 \end{array}$
Voting on 3 LVLMs	83.2 ± 1.2	83.0 ± 1.0	48.6 ± 2.8

observed in the image, avoiding assumptions about materials, functions, or contexts. If there are
any uncertainties about what an object is, describe its visual characteristics (e.g., 'a circular object
with a smooth surface') without inferring its purpose or identity. Avoid creative or hypothetical
descriptions, and focus on observable details only.

1152 With two different settings:

• **Direct Prompting**: The original input query was replaced with the prompts as described.

Prompts as Additional Guidance: We incorporated the prompt as supplemental context to guide the models in generating outputs.

1156

1134

The results demonstrate that incorporating the prompt can effectively enhance the recall performance for some models (e.g., LLaVA-v1.5 improves from 41.1 to 52.3 in recall). However, it did not consistently reduce hallucinations across all metrics and the performance on CHAIR scores (e.g., C_s , C_i) dropped. Meanwhile, MARINE significantly outperforms the prompting baseline approaches on CHAIR. We note the following differences between our method and prompting method:

Prompting method depends heavily on the instruction-following ability of the model. While it might mitigate the hallucination to a mild extent for strong models (e.g., LLaVA-v1.5), it may cause a weak model to hallucinate even more (e.g., LLaVA). Models also require more sophisticated fine-tuning approaches to generate better and more precise response conditioned on the prompts, as discussed in Deng et al. (2024). In contrast, our method directly addresses deficiencies in the model's vision capabilities by introducing stronger vision guidance. This makes our approach more effective even for weaker models and more cost-efficient.

Unlike prompting methods, which need to be tailored to specific tasks or datasets, our method generalizes effectively across models and datasets, reducing hallucinations while maintaining competitive recall.

1171

1173

1172 D.1.2 MAJORITY VOTING AMONG LVLMS

Ensembling different LVLMs is a strong baseline, as majority voting among models enhances 1174 robustness. However, it is less frequently used in practice due to the significant cost of acquiring 1175 and loading multiple LVLMs, which is substantially higher than that of additional vision models. 1176 In Table 19, we show the POPE results of the average performances of ensembling all poosible 1177 combinations of 2 and 3 LVLMs among the 5 LVLMs that we experimented with. Notably, our 1178 method outperforms the ensemble of two LVLMs. While the ensemble of three LVLMs achieves 1179 higher scores, it comes with significantly higher computational and memory costs, making it much 1180 less practical for many real-world scenarios. Meanwhile, our method requires only 30% more 1181 GPU memory than the plain LVLM during inference. For large-batch inference and online chatbot deployment, ensembling multiple LVLMs is much less feasible due to the substantial increase in 1182 memory consumption. In contrast, MARINE remains an accessible, efficient, and effective approach. 1183 1184 Furthermore, since most LVLMs use the same CLIP model as their vision component, ensembling 1185 multiple LVLMs primarily combines the language outputs, aiming for consistency across different LLMs—a strategy proven effective in the textual domain. This makes the ensemble of different 1186 LVLMs complementary to our method, which focuses on enhancing the vision component by 1187 incorporating multiple vision models.

Model	w/MARINE	Bleu-1 ↑	Bleu-2	↑ Bleu	-3 ↑	Bleu-4 ↑	ROUG	$E_L \uparrow$	CIDEr ↑	Average ↑
GIT-Large-COCO		35.68	23.01	15.	52	10.85	33.4	3	1.12	19.94
LLaVA	×	14.06	7.12	3.	2	1.90	22.0	6	0.08	8.16
LLaVA mPLUG-Owl2	√ X	18.59	9.96 25.16	5.4	F7 57	3.04	26.0	2 6	0.21	21.70
mPLUG-Owl2	√ √	39.51	24.37	15.	93	10.70	36.0	1	1.03	21.26
Table 21: Exper	iments on dyi	namic guio	dance st	rength	based	on confi	dence s	cores	on CHA	IR metrics
—	Method			LLaV	4	r	nPLUG-	Owl2		
_	CHAIR		$C_s \downarrow$	$C_i\downarrow$	Recal	$1\uparrow C_s\downarrow$	$C_i\downarrow$	Reca	11↑	
	Fix Guidance St Dynamic Guida	trength nce Strength	17.8 1 14.8	7.2 6.5	50.8 49.9	4.2 5.0	2.3 2.6	41. 41.	4 0	
– Table 22: Exper	iments on dy	namic gui	dance s	trength	based	l on conf	idence	scores	s on POI	PE metrics
Method			L	LaVA			mPI	LUG-O	wl2	
POPE		Accu	racy ↑	F1 ↑	Yes Ra	tio Acc	uracy ↑	F1↑	Yes Rat	io
Fix Guid Dynamic	ance Strength Guidance Strer	60 ngth 71	6.9 .97	72.9 74.48	72.3 59.8	3 8	32.8 3 3.3	82.7 83.2	49.2 49.4	
this approach to challenging. De remains an unex question format	a diverse rangeveloping a fe pored area f s and tasks, d	ge of tasks easible an or future r lemonstrat	—inclu d effect research ing its	ding ins ive me i. In con versatil	truction thod f ntrast, ity and	on-follov or ensen MARIN d practica	ving tas ibling E easil ality.	sks—is LVLN y geno	s signific Is in suc eralizes t	antly more h contexts to different
D.1.3 DISUC	SSION ON SP	PECIALIZE	ed Moi	DELS FO	or Im	AGE CA	PTION	NG		
notably outperfo the time, achiev Table 20.	orms many tra ed state-of-th	aditional in e-art resul	mage ca lts with	aptionin a signi	g met ficant	hods acr margin.	oss mu The re	ltiple sults a	benchma ire demo	irks and, at nstrated in
We make the fol	llowing key n	otes:								
• While cialized 19.94 a the rap task-sp	GIT-Large-C d training, ne werage score id advanceme ecific models	COCO outpower LVLM) even cor ent of LVI S.	perform Ms like npared LMs (ba	ns the c mPLU to spec ased on	lder I G-Ow ialize powe	LaVA n 12 achiev d captior ful LLM	nodel a ve bette ning mo backb	s expe er perf odels. one L	ected giv formance This der lama) in	ven its spe- e (21.70 vs monstrates exceeding
• As we serves mainta model,	previously p to demonstra ins LVLMs' p our method	pointed ou te that our performance even furth	it in the r metho ce on bi er impr	e paper od, whil roader t oves Ll	(line e prin asks v LaVA'	414-416 narily tar vithout si s perform), imag geting ignifica nance a	ge cap halluc nt trac icross	otioning cination le-offs. I all metri	evaluation mitigation For LLaVA ics.
For LL MARI model.	aVA, MARI NE, performa	NE impro ince remai	ves per ins large	forman ely stab	ce acr le and	oss all m l continu	etrics. es to ou	For m atperfo	PLUG- orm the s	Owl2 with specialized
The competitive supports our foc	e performanc cus on enhanc	ce of LVL cing these	Ms, co models	ombine rather	d with than s	n their fl pecialize	exibilit d ones	ty acr	oss mult	tiple tasks
D.2 Dynami	C GUIDANCE	e Streng	TH							
We conducted a evaluating both	dditional exp CHAIR and 1	periments POPE met	on dyn trics, as	amic g shown	uidan in Tal	ce streng ble 21 an	gth bas d 22. S	ed on Specifi	confider cally:	nce scores
• Fix Gu tion red	idance Stre duction and in	ngth uses a nstruction	a fixed g s adhere	guidanc ence.	e strei	ngth of 0	.7, sele	cted to	o balance	e hallucina
• Dynan mean c	nic Guidance confidence sco	e Strengtl ore (s) of t	h adjust the imag	ts the g ge-grou	uidano nding	ce streng models	th dyna to a rar	amica ige of	lly by m (0.4, 0.8	apping the) using the

Table 23: Inference Latency (ms/token) Comparison. We report both the latency and the ratio to the latency of greedy decoding of the original LVLM model.

	Greedy	LURE	Woodpecker*	VCD	OPERA	Offline MARINE	Online MARINE
Training Cost	0	10min on A100 80G	0	0	0	0	0
nference Latency	26.3 (×1.0)	179.9 (×6.84)	94.5 (×3.59)*	53.4 (×2.03)	185.1 (×7.0)	52.2 (×1.98)	52.23 (×1.985)

1248Table 24: Peak GPU Memory Usage during Inference (GB) of MARINE compared to greedy decoding1249and VCD.

Metric	Greedy	VCD	MARINE (Ours)
Peak GPU Memory Usage	23.53	20.73 (×0.88)	30.78 (×1.30)

formula

1254

1255 1256 $\gamma' = 0.4 + \frac{(0.8 - 0.4) \cdot (s - s_{\min})}{s_{\max} - s_{\min}}.$

A higher confidence score indicates more reliable guidance and corresponds to a stronger guidance strength. The results show that dynamic guidance improves performance for the weaker model LLaVA, which is more susceptible to noisy guidance. For stronger models like mPLUG-Owl2, a fixed guidance strength sufficiently mitigates object hallucinations with a robust performance.

1261 D.3 COMPREHENSIVE LATENCY AND MEMORY ANALYSIS

Visual prompts in MARINE are pre-generated prior to inference and were not included in the total latency reported in the original Table 5, which reflects an offline setting. We further add the online setting that accounts for the time required to process images with different visual encoders at inference time. Table 23 shows the updated latency comparison, which includes the latency for generating visual prompts. These results demonstrate that processing images adds only a negligible overhead to the overall latency.

We further measured the peak GPU memory usage throughout the inference for 500 image captioning questions using LLaVA model with batch size =16, max generation length = 64 tokens. The results are presented in Table 24. The result shows that during inference, the GPU memory usage increases by approximately 30% instead of doubling. Although we introduced additional vision models, they are relatively small compared to the large LLM backbone, resulting in only a modest increase in memory usage. Lastly, we'd highlight that our method is training-free, and thus we do not require additional memory or computation use for training, which would be significantly more demanding due to the computation of gradients.

1276 D.4 FURTHER STUDY ON GUIDANCE STRENGTH

Figure 8 illustrates how varying guidance strength affects the quality of LLaVA's output text in both the LLaVA-QA90 task and the image captioning task (maximum generation length = 256). Our experiments demonstrate that a guidance strength of 1 does not yield the best image captioning performance (i.e., quality of the generated text). In LLaVA-QA90 task, a guidance strength γ between 0.5 and 0.7 provides the most enhanced text quality. This aligns with a common concern in methods employing classifier-free guidance, where excessively strong guidance can divert the generation process.

Additionally, we conducted experiments comparing the overall generation quality of LLaVA using GPT-4V as a judge, scoring outputs on a scale of 10 for accuracy and detail, which further supports our findings. The comparison between models with and without balancing the original LVLM branch is summarized in Table 25.

- Figure 9 demonstrates how overly strong guidance can reduce instruction adherence by inducing the model to include unnecessary details from the image.
- 1290
- 1291
- 1293
- 1294
- 1295



In Figure 10, we illustrate a specific example that shows how MARINE influences the logit distribution of LVLMs during text generation. Specifically, MARINE is observed to selectively target the potential hallucinated tokens, reducing their original probabilities to mitigate the risk of hallucination in the generated text. For instance, in the provided example, the probability of "fork" is significantly lowered with MARINE, which would have originally resulted in a hallucinated object. Conversely, standard language elements such as "various", an adjective describing the overall image context, and Table 25: Results of GPT-4V-aided evaluation. The accuracy and detailedness metrics are on a scale of 10, and a higher score indicates better performance. The symbols \times and \checkmark indicate performance metrics without and with our method, respectively.

Tack	Matrice	LLaVA			
Task	wieutes	$\mathbf{X}(\gamma = 1)$	$\checkmark(\gamma=0.7)$		
	$\operatorname{Acc}\uparrow$	5.52	5.79		
LLavA-QA90	Detail \uparrow	4.58	4.77		
Imaga Captioning	Acc \uparrow	6.06	6.22		
image Captioning	Detail \uparrow	5.00	5.24		

"with", a crucial preposition, maintain their original probabilities. This selective nature of modulation
by MARINE ensures coherent and contextually relevant text generation that adheres to the instruction
while effectively reducing hallucinations.



(a) An example of image description where the original LLaVA outputs a hallucinated object, "fork".





(c) Probabilities of non-hallucinated words remain the same, highlighting MARINE's ability to preserve normal outputs.

Figure 10: This sample shows how MARINE controls logit distributions to mitigate hallucinations like "fork" while preserving the probabilities of "with", "various" during generation.

1386 1387 1388

1389

1384 1385

1372

1373

1374 1375

1380 1381 1382

E.2 DISCUSSION ON FINE-TUNING METHODS.

1390 The examples depicted in Figure 11 illustrate that LURE, at times, fails to adhere to the given 1391 instructions when correcting LVLM generations. Despite receiving concise image descriptions generated based on instructions for short responses, LURE predominantly overwrites them with 1392 excessively long responses that contain information irrelevant to the instruction. Furthermore, LURE 1393 fails to adequately address the binary question format of POPE, as LURE fixates on extended 1394 descriptions without responding with "yes" or "no", making its evaluation using POPE impractical. 1395 This issue can be prevalent in small-scale fine-tuning methods, where the limited variety of the 1396 specifically tailored fine-tuning dataset harms the model's performance on other tasks. In contrast, the training-free approach of MARINE demonstrates effective mitigation of hallucinations across a 1398 variety of question formats. 1399

1400 E.3 EXTENDED ANALYSIS IN ABLATION STUDY 1401

Additional experimental results explore the noise intensity of object grounding features, which are
 examined across LLaVA, InstructBLIP, and mPLUG-Owl2, with findings presented in Figures 12, 14, and 15.



Figure 11: Example responses to an image-question pair. The LURE-corrected output deviates from the original question, offering irrelevant descriptions without directly addressing the query. Woodpecker hallucinates the existence of two beds while there is only one bed in the figure. In contrast, MARINE maintains the original answer's style and adheres to the user's instruction while eliminating hallucination.

1421 This variation is achieved by implementing four confidence thresholds (0.5, 0.7, 0.9, and 0.95) in the DETR model predictions (with MARINE-Truth serving as an ideal reference), where higher 1422 thresholds correspond to lesser, yet higher-quality, visual information. Our findings highlight two 1423 significant insights. Firstly, an increase in the quality of visual information correlates with a noticeable 1424 decrease in hallucinations produced by the LVLMs. A lower threshold, which allows for more visual 1425 information but also includes noisier content, could potentially result in an increased occurrence of 1426 hallucinations. Furthermore, lower-quality visual information is associated with enhanced Recall. 1427 This suggests that LVLMs under guidance, despite the presence of noisy visual inputs, tend to focus 1428 more on the visual details (i.e., objects), resulting in more elaborate descriptions.



Figure 12: LLaVA's performance on CHAIR according to different noise intensity of object grounding features in MARINE. We consider four confidence thresholds (0.5, 0.7, 0.9, and 0.95) for DETR to vary the noise intensity.



Figure 13: An example of the negative impact of excessive guidance on LVLM's ability to follow instructions accurately. While the response with $\gamma = 1$ identifies more existing objects, it introduces irrelevant information to the instruction.

1457

1441

1442 1443

1444

1445

1446 1447

1448

1449

1450

1451





