

HALC: OBJECT HALLUCINATION REDUCTION VIA ADAPTIVE FOCAL-CONTRAST DECODING

Zhaorun Chen^{1,*} Zhuokai Zhao^{1,*} Hongyin Luo² Huaxiu Yao³ Bo Li^{1,4} Jiawei Zhou⁵

¹University of Chicago ²Massachusetts Institute of Technology ³UNC-Chapel Hill

⁴University of Illinois at Urbana-Champaign ⁵Toyota Technological Institute at Chicago

{zhaorun, zhuokai, bol}@uchicago.edu, jzhou@ttic.edu

ABSTRACT

While large vision-language models (LVLMs) have demonstrated impressive capabilities in interpreting multi-modal contexts, they invariably suffer from object hallucinations (OH). We introduce **HALC**, a novel decoding algorithm designed to mitigate OH in LVLMs. HALC leverages distinct fine-grained optimal visual information in vision-language tasks and operates on both local and global contexts simultaneously. Specifically, HALC integrates a robust auto-focal grounding mechanism (locally) to correct hallucinated tokens on the fly, and a specialized beam search algorithm (globally) to significantly reduce OH while preserving text generation quality. Additionally, HALC can be integrated into any LVLMs as a plug-and-play module without extra training. Extensive experimental studies demonstrate HALC’s effectiveness in reducing OH, outperforming state-of-the-arts across four benchmarks. Code is released at <https://github.com/BillChan226/HALC>.

1 INTRODUCTION

The confluence of natural language processing (NLP) and computer vision (CV) has undergone a transformative shift over the past years with the introduction of vision-language models (VLMs) (Long et al., 2022; Zhu et al., 2023; Liu et al., 2023b). Although VLMs have shown exceptional proficiency in integrating and interpreting intricate data across both textual and visual modalities, a significant challenge emerged as the phenomenon of *object hallucination (OH)*, where VLMs erroneously generate hallucinated objects and descriptions within their outputs (Rohrbach et al., 2018). Based on the different parts of the sentences that are being hallucinated, OH can be categorized into three types: *object existence*, *attribute*, and *relationship* hallucinations (Gunjal et al., 2023; Zhai et al., 2023).

While OH can be attributed to various factors (e.g. inherent biases related to co-occurrence (Biten et al., 2022; Zhou et al., 2023), visual uncertainty (Leng et al., 2023)) and exhibits certain patterns (e.g. knowledge aggregation (Huang et al., 2023), post-positioned (Zhou et al., 2023)), we conclude that its fundamental cause is the autoregressive nature of VLMs generation, where they increasingly rely on textual information while unavoidably reducing reliance on the visual input. This is especially obvious when longer responses are generated, which explains the correlation between higher OH and larger token lengths (Huang et al., 2023). A detailed literature review can be found in Appendix A.

To mitigate the disproportionate reliance on the textual and visual information during the autoregressive text generation, the process can be enhanced by continuously incorporating targeted visual information. As faithful text generations should guarantee that object-related text tokens are well grounded in the visual input, we hypothesize that the generation can benefit from focusing more on the fine-grained visual context for different object-related tokens.

To this end, we introduce Object **H**allucination Reduction through **A**daptive **F**ocal-**C**ontrast decoding, **HALC**, a novel decoding strategy designed to effectively counter OH and can be easily integrated into any open-source LVLMs such as MiniGPT-4 (Zhu et al., 2023; Chen et al., 2023), LLaVA (Liu et al., 2023b) and mPLUG-Owl2 (Ye et al., 2023). Specifically, HALC operates by identifying a token-wise *optimal visual context* to provide the most informative visual grounding while decoding a specific token. Consequently, HALC can uniquely addresses all three types of OH (existence, attribute, and relationship) while preserving linguistic quality in both local and global levels; locally,

*Equal contribution.

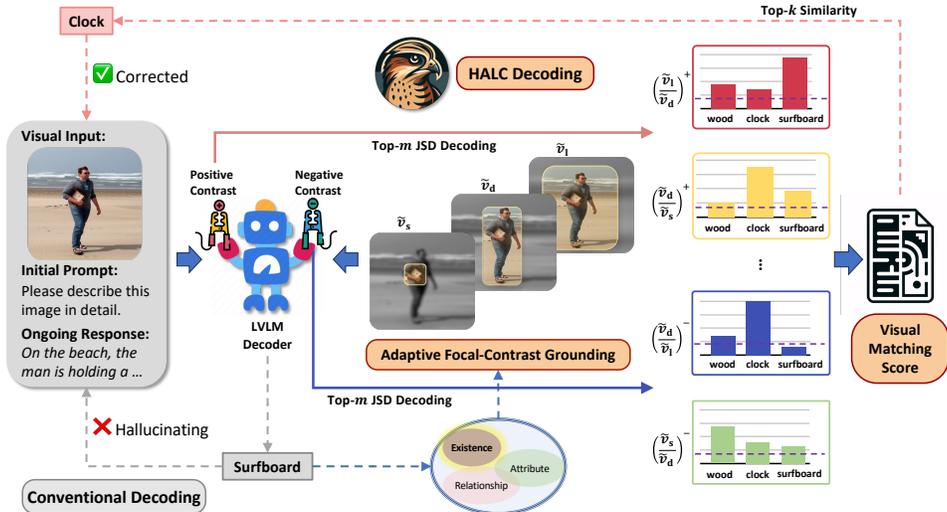


Figure 2: HALC improves text generation from images (e.g., a man holding a clock on the beach) by addressing potential errors, such as mistaking a “clock” for a “surfboard.” HALC searches for the optimal visual context for a token by first identifying its visual grounding, then sampling multiple related FOVs to obtain their logits distributions. The optimal logits distributions are approximated by the largest JSD between contrasted pairs.

it employs an *adaptive focal-contrast grounding* mechanism to locate the fine-grained optimal visual information to correct each generated token that might be hallucinating; and globally, it incorporates a *matching-based beam search* that utilizes a visual matching score to steer the generation of the final outputs to balance both OH mitigation and text generation quality.

2 OH AND FINE-GRAINED VISUAL KNOWLEDGE

Problem Formulation. For an image-grounded text generation task, a θ -parameterized LVLM $\mathcal{M}_\theta^{\text{LVLM}}$ often generate texts in an auto-regressive manner. Given a textual query x and an input image v , v is first processed by a vision encoder into a visual embedding, then transformed by a multi-modal projector together with the query x , and finally decoded into a textual response y . OH happens when some parts of the text generation y is inconsistent with the input image v . The goal of HALC is to minimize the occurrence of OH tokens and preserve the faithfulness to v in x , while maintaining a high-quality generation of text y . A detailed problem formulation can be found in Appendix B.

Fine-grained Visual Knowledge. To mitigate the disproportionate reliance on the textual and visual information during the autoregressive text generation, the process can be enhanced by continuously incorporating targeted visual information. As faithful text generations should guarantee that object-related text tokens are well grounded in the visual input, we hypothesize that the generation can benefit from focusing more on the *fine-grained visual context* for different object-related tokens. We verify our hypothesis through an empirical pilot study: Fig. 1 shows OH percentages when we feed the greedy decoding with or without brute-force searched optimal visual contexts on the OH subset of the MME benchmark (Fu et al., 2023)(the implementation details can be found in Appendix G). We can see that incorporating such optimal visual contexts can eliminate over 84.5% of the OH. This observation leads to the key insight in HALC that mitigating OH lies in identifying a token-wise optimal visual context to provide the most informative visual grounding while decoding a specific token, which is achieved through its *adaptive focal-contrast decoding* module.

3 HALC

Deriving from the above statistical analysis of the effectiveness of fine-grained visual context in correcting OH, we propose **HALC**. A schematic overview of HALC is shown in Fig. 2. HALC operates at the token level during generation, with reliance on fine-grained visual information

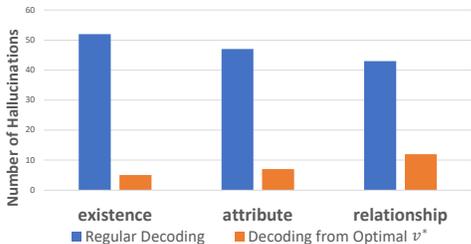


Figure 1: On average, over 84.5% of the OH are reduced by leveraging some optimal visual context v^* .

represented by samples of different visual contexts. By recomputing the token distributions from different visual context inputs and contrasting them, object-related token probabilities are redistributed to reduce hallucinations dynamically within the generation steps.

Object-related token identification. To efficiently target likely sources of OH, we initially pinpoint tokens linked to objects for HALC processing. Specifically, at each generation step, we determine the part-of-speech (POS) tag of the current token from the model $\mathcal{M}_\theta^{\text{LVLML}}$. If the token is a noun, adjective, adverb, number, verb, pronoun, or preposition—indicating potential for object, attribute, or relationship hallucinations—we re-generate the token with HALC. For instance, as in Fig. 2, the token “surfboard” might be flagged for potential object existence hallucination. Notice that we base the decision to reprocess a token on its syntactic category, without assuming it’s already hallucinating..

Visual context retrieval. To obtain detailed visual information for a token, we identify a corresponding visual context window $v_d = (w_d, h_d, p_d)$ defined by its width, height, and center point. We use a zero-shot detector, such as Grounding DINO (Liu et al., 2023c) or OWLv2 (Minderer et al., 2023), to pinpoint the token’s location in the image. While these detectors are primarily used for object detection, they’re also capable of providing visual references for adjectives or prepositional phrases. This is due to their pre-training goal of linking text descriptions to image regions, encompassing attributes and relationships in addition to object identification (Liu et al., 2023c).

Adaptive focal-contrast grounding. While off-the-shelf detectors establish a meaningful reference v_d within the original image input v , it is often not the optimal visual context for decoding. As shown in Fig. 3, we demonstrate that the likelihood of different objects’ tokens can vary significantly across various visual context windows fed into $\mathcal{M}_\theta^{\text{LVLML}}$. For instance, in a case where the correct token “clock” is mistakenly decoded as “surfboard,” we find that an alternative visual context, labeled as v_1 , more accurately corrects this error by significantly increasing the “clock” token’s probability. This standout pattern of the “clock” token, compared to the more uniform patterns of other tokens across different visual contexts, underlines our strategy of *focal-contrast grounding*. This method aims to fine-tune the probabilities of object-related tokens by exploring and selecting from a spectrum of FOVs that sharply contrast in their decoding probabilities, thereby closely approximating the ideal visual contexts.

FOV sampling. We first sample a sequence of n FOVs, v_1, v_2, \dots, v_n , based on the initial visual context v_d . Various methods can generate these FOVs based on v_d . To attain a larger coverage of the input image quickly, one strategy of FOVs sampling is through an exponential expanding function, by setting $v_i = (w_i, h_i, p_i) = ((1 + \lambda)^i w_d, (1 + \lambda)^i h_d, p_d)$, where w_i, h_i, p_i are the width, height, and center of the FOV v_i .

Dynamic visual context selection. Based on the observation from Fig. 3, we now select a set of FOVs based on a contrastive criterion in the text decoding space to better approximate the optimal visual context for the current token. In particular, after obtaining n different FOVs, we feed these visual contexts back into the model $\mathcal{M}_\theta^{\text{LVLML}}$, resulting in n different probability distributions $p_i = p_\theta(\cdot|v_i, x, y_{<t})$ with $i = 1, 2, \dots, n$. Between any two candidate FOVs, we adopt the following distance measure for the discrepancy between their decoded token probability distributions

$$d(v_i, v_j) = \text{JSD}(p_\theta(\cdot|v_i, x, y_{<t}) \parallel p_\theta(\cdot|v_j, x, y_{<t})) \quad (1)$$

where JSD is the Jensen-Shannon divergence, a symmetric metric that measures the difference between two distributions. With the idea that more different FOV pairs are more likely to include the optimal visual context for the current victim token generation, we dynamically select the top m pairs with the largest distance according to Eq. (1).

Contrastive decoding. After obtaining top m visual context pairs with most discrepancies in influencing the token output, we contrast the decoding probability distributions (p_i, p_j) within each



Figure 3: In the FOV space, clear objects (“beach”, “man”) maintain stable, high likelihoods, while hallucinated objects, like “book” or “surfboard,” show erratic or shifting likelihoods. Incorrectly generated tokens, or “victim tokens” (e.g., “clock”), typically exhibit a sharp peak in likelihood, signaling a local maximum.

pair in order to amplify the information residing in one visual context over the other. This would potentially recover the victim token over the hallucinated token as the victim token enjoys a sharper contrast in the probability comparisons, especially when one of the visual contexts under comparison is near the optimal grounding. Specifically, we redistribute the probabilities based on the contrast in log space (Li et al., 2022b) for a given FOV pair (v_i, v_j) , resulting in $p_{v_i/v_j}(\cdot|v_i, v_j, x, y_{<t}) \propto \exp \left[(1 + \alpha) f_{\theta}(\cdot|v_i, x, y_{<t}) - \alpha f_{\theta}(\cdot|v_j, x, y_{<t}) \right]$ where f_{θ} again is the logit distribution, α is the amplification factor where larger α indicates a stronger amplification of the differences.

Unlike traditional uni-modal contrastive decoding methods (Chuang et al., 2023; Gera et al., 2023; Shi et al., 2023) that distinguish between expert and amateur distributions based on the assumption that the final or context-aware layer has more accurate knowledge, our approach to determining an expert distribution among FOV pairs is complex due to the optimal visual context often being between expanding FOVs. This can lead to OH with too much or too little context. Without knowing the exact location of the optimal context, we contrast each FOV pair bi-directionally, incorporating both positive (larger over smaller FOV) and negative (smaller over larger FOV) contrasts to ensure complete FOV representation. This yields $2m$ candidate tokens from individual decodings, later refined by a matching-based beam search algorithm.

Matching-based beam search. Our adaptive focal-contrast grounding corrects individual tokens during generation, while a sequence-level beam search with a beam size of k ensures overall text quality. At each HALC decoding step, k beam sequences produce $2mk$ token candidates from the top m focal-contrast pairs. Unlike traditional beam search methods that use only textual information, we select the top k beams from $2mk$ candidates based on a global visual matching score, comparing text sequence similarity to the original image, ensuring diversity and accuracy in the generated text. The BLIP model (Li et al., 2022a) is used for encoding both text and image to compute similarity scores. The complete HALC process is detailed in the Appendix C. We conduct both theoretical certified robustness analysis and empirical analysis of our optimal visual contexts approximation, which is the most important component of HALC, in Appendix D and Appendix G respectively.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Benchmarks. We evaluate HALC on three benchmarks including (1) quantitative metrics CHAIR (Rohrbach et al., 2018) and POPE (Li et al., 2023) on MSCOCO (Lin et al., 2014) dataset; (2) general-purposed Multimodal Large Language Model Evaluation (MME) (Fu et al., 2023) benchmark; and (3) qualitative evaluation benchmark LLaVA-Bench (Liu et al., 2023a). These experiments comprehensively assess HALC’s capability on reducing OH in image captioning, visual-question answering (VQA) and more challenging tasks that generalize to novel domains.

Baselines. To effectively evaluate HALC, besides regular greedy decoding and beam search baselines, we further involve layer-wise contrastive decoding SOTA DoLa (Chuang et al., 2023), as well as SOTA methods specifically designed to mitigate OH, including OPERA (Huang et al., 2023), VCD (Leng et al., 2023), Woodpecker (Yin et al., 2023) and LURE (Zhou et al., 2023) in our analysis. All the results are acquired and benchmarked consistently within our unified implementation. Please refer to Appendix F for the detailed setting of our experiments including hyper-parameters.

LVLm Backbones. Three LVLms (MiniGPT-4 V2 (Chen et al., 2023), LLaVA-1.5 (Liu et al., 2023b), mPLUG-Owl2 (Ye et al., 2023)) are assessed for both HALC and all above baselines except Woodpecker and LURE, where Woodpecker utilizes ChatGPT (Brown et al., 2020) during its self-correction process and LURE distills an extra reviser model from GPT-4 (Achiam et al., 2023).

4.2 RESULTS

Following existing evaluation procedures (Huang et al., 2023; Yin et al., 2023; Liu et al., 2023b), we randomly sampled 500 images from the validation split of MSCOCO (Lin et al., 2014) and conduct evaluations with both CHAIR and POPE. For each metric, we repeat the experiments five times with different random seeds and report average and standard deviations of all the runs.

CHAIR. The Caption Hallucination Assessment with Image Relevance (CHAIR) (Rohrbach et al., 2018) evaluates the occurrence of OH in image captioning tasks. It measures the extent of OH by determining the proportion of mentioned objects that are absent in the actual label set. CHAIR

includes two metrics: CHAIR_S, assessing sentence-level hallucinations, and CHAIR_I, assessing object instance-level hallucinations. Lower scores in either metric indicate fewer hallucinations. Besides CHAIR_S and CHAIR_I, we also report BLEU (Papineni et al., 2002) as an assessment of the text generation quality. Table 1 demonstrates that our proposed HALC consistently outperforms all the existing methods by a large margin. Notably, a major advantage of HALC is its strong robustness, as can be observed by its much lower standard deviations, especially when compared to the non-OH specific baselines.

Table 1: CHAIR evaluation results on MSCOCO dataset of LVLMs with different decoding baselines and SOTAs designed for mitigating OH. Lower CHAIR_S and CHAIR_I indicate less OH. Higher BLEU generally represent higher captioning quality, although existing work has reported weak correlation between CHAIR and text overlapping quality metrics. Bold indicates the best results of all methods.

Method	MiniGPT-4			LLaVA-1.5			mPLUG-Owl2		
	CHAIR _S ↓	CHAIR _I ↓	BLEU ↑	CHAIR _S ↓	CHAIR _I ↓	BLEU ↑	CHAIR _S ↓	CHAIR _I ↓	BLEU ↑
Greedy	30.87 \pm 5.45	12.33 \pm 2.07	14.33 \pm 0.00	20.80 \pm 0.08	6.77 \pm 0.07	15.93 \pm 0.00	23.20 \pm 0.35	8.33 \pm 0.28	15.37 \pm 0.00
Beam Search	29.56 \pm 6.09	11.36 \pm 0.99	14.94 \pm 0.00	18.67 \pm 0.38	6.30 \pm 0.05	16.17 \pm 0.00	21.67 \pm 1.61	7.63 \pm 0.40	15.77 \pm 0.00
DoLA	30.87 \pm 2.52	11.70 \pm 0.13	14.93 \pm 0.00	21.00 \pm 0.67	6.70 \pm 0.38	15.93 \pm 0.00	24.60 \pm 0.24	8.73 \pm 0.30	15.40 \pm 0.00
OPERA	30.00 \pm 0.43	11.67 \pm 0.22	14.87 \pm 0.00	21.13 \pm 0.12	6.73 \pm 0.18	16.27 \pm 0.01	22.13 \pm 0.86	7.57 \pm 0.16	15.53 \pm 0.00
VCD	30.27 \pm 0.44	12.60 \pm 0.45	14.33 \pm 0.00	23.33 \pm 5.66	7.90 \pm 0.53	14.67 \pm 0.01	27.27 \pm 7.32	9.73 \pm 1.22	14.40 \pm 0.00
Woodpecker	28.87 \pm 2.20	10.20 \pm 0.85	15.30 \pm 0.01	23.85 \pm 4.62	7.50 \pm 0.01	17.05 \pm 0.00	26.33 \pm 1.98	8.43 \pm 0.80	16.43 \pm 0.00
LURE	27.88 \pm 2.25	10.20 \pm 0.85	15.03 \pm 0.11	19.48 \pm 2.35	6.5 \pm 0.38	15.97 \pm 0.01	21.27 \pm 0.06	7.67 \pm 0.16	15.65 \pm 0.05
HALC	17.80 \pm 0.03	8.10 \pm 0.14	14.91 \pm 0.00	13.80 \pm 0.08	5.50 \pm 0.14	16.10 \pm 0.01	17.33 \pm 4.30	7.43 \pm 0.11	16.27 \pm 0.00

POPE. Polling-based Object Probing Evaluation (POPE) (Li et al., 2023) evaluates OH via a streamlined approach, which incorporates a list of yes-or-no questions to prompt LVLMs for presence of positive and negative objects. Unlike CHAIR, POPE directly interacts with the examined large vocabulary language model (LVM), which is suitable for decoding-based baselines but less adaptable for post-hoc methods like LURE (Zhou et al., 2023). This approach also leads to greater instabilities with smaller language backbones such as LLaMA-7B, which has weaker chat capabilities. In response, we introduce *offline POPE (OPOPE)*, which retains POPE’s object sampling and yes/no queries but substitutes live interactions with offline checks. Specifically, instead of querying, "Is there a in the image?", OPOPE first obtains the LVM’s detailed image descriptions and then manually verifies the presence of sampled objects in these captions to compute the scores. The evaluation results incorporating OPOPE is shown in Table 2. HALC outperforms other methods in most of the settings.

Table 2: Proposed OPOPE evaluation results on MSCOCO dataset of LVLMs with different decoding baselines and SOTAs designed for mitigating OH. Higher accuracy, precision, and F score indicate better performance. Bold indicates the best results of all methods.

Method	MiniGPT-4			LLaVA-1.5			mPLUG-Owl2		
	Accuracy ↑	Precision ↑	F _{β=0.2} ↑	Accuracy ↑	Precision ↑	F _{β=0.2} ↑	Accuracy ↑	Precision ↑	F _{β=0.2} ↑
Greedy	66.78 \pm 1.27	90.43 \pm 25.1	85.79 \pm 18.7	70.56 \pm 1.51	91.08 \pm 20.6	87.72 \pm 16.3	69.77 \pm 1.18	91.07 \pm 17.8	87.45 \pm 13.9
Beam Search	67.22 \pm 0.74	91.20 \pm 14.4	86.57 \pm 10.8	69.87 \pm 1.37	91.72 \pm 20.4	88.01 \pm 15.97	69.20 \pm 0.90	91.90 \pm 15.1	87.91 \pm 11.7
DoLA	67.06 \pm 1.19	90.84 \pm 23.1	86.22 \pm 17.3	70.69 \pm 1.50	90.87 \pm 19.8	87.59 \pm 15.74	70.17 \pm 1.69	91.97 \pm 24.5	88.30 \pm 19.26
OPERA	67.26 \pm 1.04	90.76 \pm 20.0	86.25 \pm 15.0	69.73 \pm 1.34	91.10 \pm 19.4	87.46 \pm 15.3	69.26 \pm 0.45	93.06 \pm 8.01	88.83 \pm 6.14
VCD	65.78 \pm 0.96	90.02 \pm 20.7	85.00 \pm 15.1	70.67 \pm 1.22	91.62 \pm 16.7	88.19 \pm 13.3	69.81 \pm 0.65	92.70 \pm 11.0	88.76 \pm 8.49
Woodpecker	67.78 \pm 0.88	91.33 \pm 16.66	86.91 \pm 12.6	69.80 \pm 0.54	91.80 \pm 8.41	88.04 \pm 6.56	68.90 \pm 1.02	92.22 \pm 17.98	88.05 \pm 13.77
LURE	68.14 \pm 0.99	90.95 \pm 17.34	86.76 \pm 13.23	70.00 \pm 1.53	90.89 \pm 21.9	87.38 \pm 17.3	69.24 \pm 1.60	90.54 \pm 23.37	86.85 \pm 18.28
HALC	66.76 \pm 0.68	91.95 \pm 15.0	86.92 \pm 11.1	70.59 \pm 0.82	92.94 \pm 12.18	89.22 \pm 9.55	70.12 \pm 0.98	91.94 \pm 15.1	88.26 \pm 11.85

More results. More detailed results on CHAIR, POPE, and additional evaluations on LLaVA-Bench can be found in Appendix E.

5 CONCLUSION

We present HALC, a novel decoding algorithm designed to mitigate OH in LVLMs. HALC operates on both local and global levels, integrating a robust adaptive focal-contrast grounding mechanism to better utilize fine-grained visual information for correcting hallucinated tokens, and a specialized beam search algorithm that promotes further visually matched generations. Comprehensive experiments demonstrate that HALC effectively reduces OH, achieving SOTA performance while preserving sequence generation quality, and can be conveniently integrated into existing LVLMs without additional training or data. A benchmarking tool was also built to support convenient comparisons across all available OH reduction strategies comprehensively.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 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.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 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.
- Zhaorun Chen, Zhuokai Zhao, Zhihong Zhu, Ruiqi Zhang, Xiang Li, Bhiksha Raj, and Huaxiu Yao. Autoprml: Automating procedural supervision for multi-step reasoning via controllable question decomposition. *arXiv preprint arXiv:2402.11452*, 2024.
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. Dola: Decoding by contrasting layers improves factuality in large language models. *arXiv preprint arXiv:2309.03883*, 2023.
- Chenhang Cui, Yiyang Zhou, Xinyu Yang, Shirley Wu, Linjun Zhang, James Zou, and Huaxiu Yao. Holistic analysis of hallucination in gpt-4v (ision): Bias and interference challenges. *arXiv preprint arXiv:2311.03287*, 2023.
- Wenliang Dai, Zihan Liu, Ziwei Ji, Dan Su, and Pascale Fung. Plausible may not be faithful: Probing object hallucination in vision-language pre-training. *arXiv preprint arXiv:2210.07688*, 2022.
- Imant Daunhawer, Thomas M Sutter, Kieran Chin-Cheong, Emanuele Palumbo, and Julia E Vogt. On the limitations of multimodal vaes. *arXiv preprint arXiv:2110.04121*, 2021.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- Ariel Gera, Roni Friedman, Ofir Arviv, Chulaka Gunasekara, Benjamin Sznajder, Noam Slonim, and Eyal Shnarch. The benefits of bad advice: Autocontrastive decoding across model layers. *arXiv preprint arXiv:2305.01628*, 2023.
- Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, and Dinesh Manocha. Tianyi Zhou. Hallusionbench: An advanced diagnostic suite for entangled language hallucination & visual illusion in large vision-language models. *arXiv e-prints*, pp. arXiv–2310, 2023.
- Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and preventing hallucinations in large vision language models. *arXiv preprint arXiv:2308.06394*, 2023.
- Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. Opera: Alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation. *arXiv preprint arXiv:2311.17911*, 2023.
- 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.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pp. 12888–12900. PMLR, 2022a.

- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization. *arXiv preprint arXiv:2210.15097*, 2022b.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023a.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023b.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023c.
- Siqu Long, Feiqi Cao, Soyeon Caren Han, and Haiqin Yang. Vision-and-language pretrained models: A survey. *arXiv preprint arXiv:2204.07356*, 2022.
- Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. Scaling open-vocabulary object detection. *arXiv preprint arXiv:2306.09683*, 2023.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. *arXiv preprint arXiv:1809.02156*, 2018.
- Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Scott Wen-tau Yih. Trusting your evidence: Hallucinate less with context-aware decoding. *arXiv preprint arXiv:2305.14739*, 2023.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4566–4575, 2015.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pp. 38–45, 2020.
- Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. *arXiv preprint arXiv:2311.04257*, 2023.
- 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. *arXiv preprint arXiv:2310.16045*, 2023.

- Bohan Zhai, Shijia Yang, Chenfeng Xu, Sheng Shen, Kurt Keutzer, and Manling Li. Halle-switch: Controlling object hallucination in large vision language models. *arXiv e-prints*, pp. arXiv-2310, 2023.
- Yiming Zhang, Zhuokai Zhao, Zhaorun Chen, Zhili Feng, Zenghui Ding, and Yining Sun. Rankclip: Ranking-consistent language-image pretraining. *arXiv preprint arXiv:2404.09387*, 2024.
- Zhuokai Zhao, Harish Palani, Tianyi Liu, Lena Evans, and Ruth Toner. Multi-modality guidance network for missing modality inference. *arXiv preprint arXiv:2309.03452*, 2023.
- 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. *arXiv preprint arXiv:2310.00754*, 2023.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.

A RELATED WORK

Object hallucination (OH). OH refers to the phenomenon where vision-language models (VLMs), including both the earlier BERT-based models (Li et al., 2019; Radford et al., 2021) and the more recent LVLMs (Liu et al., 2023b; Zhu et al., 2023), erroneously generate unfaithful contents. More specifically, Gunjal et al. (2023) and Zhai et al. (2023) proposed that OH could be categorized into three types: object *existence* hallucination for the creation of non-existent objects, object *attribute* hallucination for providing misleading descriptions, and object *relationship* hallucination for depicting incorrect inter-object relationships.

Why does OH occur? OH in VLMs can be attributed to various factors, including but not limited to the inherent biases in the training data caused by co-occurrence (Biten et al., 2022; Zhou et al., 2023), visual uncertainty due to model’s statistical bias and priors (Leng et al., 2023), as well as the limitations in current models’ ability to discern context and fact accurately during the entire output generation process (Daunhawer et al., 2021). Studies have also shown that OH is not random but exhibits certain patterns and dependencies, such as its co-existence with knowledge aggregation pattern (Huang et al., 2023), and the tendency to occur with objects positioned later in the generated descriptions (Zhou et al., 2023).

A closer examination of these analysis suggests that the autoregressive nature of the LVLMs may be a fundamental factor contributing to their hallucinatory behaviors. Specifically, autoregressive decoding makes LVLMs progressively rely more on textual information including both the query x and the increasing history generations $y_{<t}$, while unavoidably reducing reliance on the visual input. This imbalance results in a significant deviation from accurate representation of the visual input, ultimately culminating in OH with behaviors and patterns observed in the aforementioned studies (Zhou et al., 2023; Leng et al., 2023). This is especially obvious when longer responses are generated, which explains the correlation between higher OH and larger maximum token lengths, as seen in (Huang et al., 2023).

OH assessment. The most well-adopted metric specifically designed to evaluate OH is CHAIR (Rohrbach et al., 2018), which was motivated after Rohrbach et al. (2018) discovered that existing metrics that measure the output’s text quality, such as CIDEr (Vedantam et al., 2015), is misleading at representing hallucinations (higher CIDEr score may correlate with higher OH). Another notable and more recent metric is POPE (Li et al., 2023), which transforms the assessment of OH into a binary classification problem where metrics such as precision, recall and accuracy are used to represent the level of OH. In our evaluations, we utilize CHAIR and propose a new metric based on POPE, named *OPOPE*, for thorough assessments of OH, while keeping the standard text generation quality metrics such as BLEU (Papineni et al., 2002), as an additional indicator to make sure little sacrifice in quality was made when mitigating OH.

Challenges and existing approaches. OH has been a persistent challenge since the earlier stages of the VLM development (Rohrbach et al., 2018; Cui et al., 2023). And it has been gaining increased attention, especially when recent research indicates that even the much more sophisticated and capable large vision-language models (LVLMs) are not immune to it (Dai et al., 2022; Li et al., 2023; Guan et al., 2023). Despite numerous advancements in LVLMs (Zhao et al., 2023; Chen et al., 2024; Zhang et al., 2024), none of them can produce faithful outputs without suffering from some level of OH. Various strategies have been developed to this matter. For instance, Zhou et al. (2023) and Yin et al. (2023) proposed post-hoc and self-correction pipelines, respectively. Huang et al. (2023) and Leng et al. (2023) developed decoding strategies emphasizing better prior utilization. While effective, these approaches often require powerful external LVLMs or additional data, limiting their adaptability.

Despite the efforts, these approaches are not yet fully satisfying in terms of eliminating OH. More importantly, they mainly focus on mitigating object existence hallucination, while assuming the attribute- and relationship-level hallucinations can be consequently corrected through autoregressive decoding. Furthermore, their reliance on more powerful external LVLMs (Yin et al., 2023), repeated processing (Zhou et al., 2023) or additional data (Gunjal et al., 2023) complicates their adaptations to existing LVLMs and restricts their use cases. The importance of OH reduction combined with the limitations in existing methods underscore the urgent need for developing novel approaches.

Distinct from these methods, HALC offers a novel decoding strategy that effectively reduces OH without necessitating extra LVLMs, training, or data. Integrating a novel adaptive focal-contrast grounding mechanism, HALC addresses both local and global contexts in OH reduction. Its compati-

bility with open-source LVLMs like MiniGPT-4 (Zhu et al., 2023) and LLaVA (Liu et al., 2023b) further enhances its applicability.

And as previous approaches often study the problem under different settings and metrics (Zhou et al., 2023; Yin et al., 2023; Huang et al., 2023; Leng et al., 2023), to promote the development of OH reduction in general, we implement an open-source platform which hosts both the proposed HALC and other methods, supporting various LVLM backbones and evaluation metrics.¹

B DETAILED PROBLEM FORMULATION

We consider an LVLM $\mathcal{M}_\theta^{\text{LVLM}}$ parameterized by θ , with a general architecture consisting of a vision encoder, a vision-text interface module, and a text decoder. For an image-grounded text generation task, given a textual query x and an input image v , v is first processed by the vision encoder into a visual embedding, then transformed by the interface module as the input to the text decoder together with the query x , and finally decoded into a textual response y autoregressively. Formally, we have

$$y_t \sim p_\theta(\cdot | v, x, y_{<t}) \propto \exp f_\theta(\cdot | v, x, y_{<t}) \quad (2)$$

where y_t denotes the t^{th} token, $y_{<t}$ is the token sequence generated up to time step t , and f_θ is the logit distribution (unnormalized log-probabilities) produced by $\mathcal{M}_\theta^{\text{LVLM}}$.

OH happens when some parts of the text generation y conflicts with the input image v . The goal of OH reduction is to minimize the occurrence of hallucination tokens and preserve the faithfulness to v when addressing the query x , while maintaining a high-quality generation of text y .

C HALC ALGORITHM

Algorithm 1 HALC Decoding

Require: LVLM $\mathcal{M}_\theta^{\text{LVLM}}$, text query x , image input v , grounding detector \mathcal{G}_d , FOV sample size n , beam size k , number of contrast FOV pairs m .

output Model response y_{new} .

```

1: repeat
2:   At every decoding step  $t$ :
3:   for  $b = 1$  to beam size  $k$  do
4:      $\mathcal{M}_\theta^{\text{LVLM}}$  decoding, obtain current token  $y_t^b$ 
5:     if  $y_t^b \in \{\text{existence, attribute, relationship}\}$  then
6:       Retrieve visual context  $v_d^b \leftarrow \mathcal{G}_d(y_t^b, v)$ 
7:     end if
8:     if  $v_d^b \neq \{\emptyset\}$  then
9:       Sample  $n$  FOVs  $v_1, \dots, v_n$  by expanding  $v_d^b$ 
10:    else
11:      Randomly sample  $n$  FOVs  $v_1, \dots, v_n$  from  $v$ 
12:    end if
13:    Compute pair-wise JSDs  $d(v_i, v_j), \forall i \neq j$ 
14:    Select top- $m$  candidate pairs
15:    for  $i = 1$  to  $m$  do
16:      Apply bi-directional contrast  $(p_{v_i/v_j}, p_{v_j/v_i})$ ,
17:      get a pair of redistributed logits
18:    end for ▷  $y_{\text{new}}^b$  with  $2m$  candidates obtained
19:    end for
20:    Select top  $k$  candidates by visual matching
21:    if  $v_d^b \neq \{\emptyset\}$  and  $y_{\text{new}}^b = y_t^b$  then
22:       $y_{\text{new}}^b \leftarrow [\text{IDK}]$  ▷  $y_t^b$  is hallucinating, but no correction token was found
23:    end if
24:     $y_t^b \leftarrow y_{\text{new}}^b$  ▷ Hallucinating token  $y_t^b$  corrected
25: until each beam has terminated

```

¹We make our codes public at <https://github.com/BillChan226/HALC>.

D THEORETICAL ANALYSIS ON FOV SELECTION

Based on our observation (in Fig. 1 and Fig. 3) that there exists some underlying optimal visual context v^* within the original image v that can largely reduce the object hallucination at the token level, our method aims to recover this optimal visual context v^* based on a sampling process conditioned on v_d . To do so, we first select the visual contexts, or FOVs, by taking a sequence of FOV samples starting from the initial v_d based on an off-the-shelf detector. While we cannot guarantee that the initial visual grounding v_d is sufficiently accurate to approximate v^* (and directly using v_d could result in unstable behaviors), we could effectively certify the robustness of our FOV sampling strategy in Theorem D.1. To preserve generality, consider the sampled FOVs are taken from a distribution $\pi(\cdot|v_d)$, where π can either follow normal distribution sampling around v_d , or obey an exponential expansion sampling strategy starting from v_d .

Theorem D.1. *Let $v^* = (w^*, h^*, p^*)$ be the optimal visual context. Assume there exists a tolerable neighborhood $\mathcal{B}(v^*, \epsilon) = \{\hat{v} : \|\hat{v} - v^*\| \leq \epsilon\}$ around v^* , such that decoding from visual contexts within the neighborhood is robust:*

$$D(p_\theta(\cdot|v^*), p_\theta(\cdot|\hat{v})) \leq \delta \ll 1, \forall \hat{v} \in \mathcal{B}(v^*, \epsilon) \quad (3)$$

where $D(\cdot, \cdot) \in [0, 1]$ is a symmetric discrepancy measure between two probability distributions, such as the Jensen-Shannon divergence, or the total variation distance.

Let $v_d = (w_d, h_d, p_d)$ be the initial detection and $v_d = v^* + \eta$ with perturbation η . The minimum deviation of token probabilities from the optimum with n samples v_1, v_2, \dots, v_n distributed according to $\pi(\cdot|v_d)$ is denoted as

$$h_\pi(v^*, n) = \min_{i=1, \dots, n} D(p_\theta(\cdot|v^*), p_\theta(\cdot|v_i)) \quad (4)$$

(a) For normal distribution sampling $\pi_g(\cdot|v_d) \sim \mathcal{N}(v_d, \sigma^2 I)$, the minimum deviation above is bounded as

$$h_{\pi_g}(v^*, n) \leq \delta + (1 - C_g(\epsilon, \eta; \sigma))^n \quad (5)$$

where $C_g(\epsilon, \eta; \sigma) \in (0, 1)$ is a constant depending on ϵ, η, σ , and the upper bound goes to δ when $n \rightarrow \infty$.

(b) For exponential expansion sampling $\pi_e(\cdot|v_d) \sim \mathcal{U}(r \in [r_{\min}, r_{\max}])$ with samples $v_r = ((1 + \lambda)^r w_d, (1 + \lambda)^r h_d, p_d)$ uniformly from the r -space, under the conditions (i) $|p_d - p^*| < \epsilon$ and (ii) $w_d/h_d = w^*/h^*$, the minimum deviation in Eq. equation 4 is bounded below

$$h_{\pi_e}(v^*, n) \leq \delta + (1 - C_e(\epsilon, v^*, v_d; \lambda))^n \quad (6)$$

where $C_e(\epsilon, v^*, v_d; \lambda) \in (0, 1]$ is a constant depending on $\epsilon, v^*, v_d, \lambda$, and the upper bound goes to δ when $n \rightarrow \infty$.

The proof of Theorem D.1 is detailed below. The neighborhood radius ϵ around the optimal v^* can be roughly interpreted as a valid range of optimal visual context to yield the correct prediction (e.g., $[v_1, v_2]$ in Fig. 3). Typically the detection perturbation $\|\eta\| > \epsilon$, making v_d outside of the ϵ -neighborhood of v^* . Through FOV sampling according to some $\pi(\cdot|v_d)$, the above theorem establishes a formal guarantee that at least one of the n samples achieves good approximation of the optimal v^* in the decoding probability space, as the deviation is closer to δ when n grows. The normal sampling distribution, concentrated around v_d , is preferred when v_d has minimal perturbations from v^* . And an exponential expansion sampling distribution, with a more averaged coverage of the sampling space, is preferable when less prior of the task is available. In practice of our algorithm, we take discrete integer values of r under the exponential expansion distribution for deterministic sampling with $n = 4$, acquiring good efficiency and performance.

Proof. Let $v^* = (w^*, h^*, p^*)$ be the optimal visual context, represented by a 3-tuple of its width, height, and center point. The corresponding optimal token decoding probability distribution is $p_\theta(\cdot|v^*)$, where θ denotes the parameters of the LVLm $\mathcal{M}_\theta^{\text{LVLm}}$, and we ignore the condition on the textual query x and previously generated tokens $y_{<t}$ for simplicity. We rely on a symmetric discrepancy measure $D(\cdot, \cdot) \in [0, 1]$ to compare the disparity between two probability distributions, such as the Jensen-Shannon divergence, or the total variation distance. We assume that the model

prediction is robust around v^* against small perturbations. In particular, we assume that there exists a tolerable small ϵ -neighborhood $\mathcal{B}(v^*, \epsilon) = \{\hat{v} : \|\hat{v} - v^*\| \leq \epsilon\}$ around v^* , such that

$$g(v^*, \hat{v}) = D(p_\theta(\cdot|v^*), p_\theta(\cdot|\hat{v})) \leq \delta \ll 1, \quad \forall \hat{v} \in \mathcal{B}(v^*, \epsilon) \quad (7)$$

Essentially, for any visual context window (or FOV) close enough to v^* , the output token probability disparity is tiny, which is likely to result no difference in greedy decoding.

From the FOV detector \mathcal{G}_d , the output visual context is denoted as $v_d = (w_d, h_d, p_d)$, which is in general not the optimal. We assume $v_d = v^* + \eta$ in the 3-tuple vector space, where η is the perturbation vector from the optimal. The detection perturbation is often large enough with $\|\eta\| > \epsilon$, making v_d outside of the ϵ -neighborhood of v^* .

$v_d \rightarrow v^*$: If we directly use the detector output v_d as an approximation of the optimal visual context v^* , the output distribution deviation from the optimum, measured by $g(v^*, v_d)$, is often unpredictable, when v_d does not fall in the hypothetical tolerable region $\mathcal{B}(v^*, \epsilon)$. An example can be seen as the inaccurate detection v_d in Fig. 3 results in the wrong token prediction *book*. This prompts the need for our proposed FOV sampling approach with the hope to find samples close to the optimal v^* .

$\pi(\cdot|v_d) \rightarrow v^*$: Thus we consider sampling conditioned on v_d in the FOV space to enhance the robustness of optimal visual context approximation, hoping to find some sample that is close to the optimal. To do this, we obtain an upper bound on the minimum deviation from the output distribution among a collection of FOV samples. Assume $\pi(\cdot|v_d) \in \Omega$ is an arbitrary sampling function conditional on the initial FOV detection v_d , where Ω denotes the sampling space over all potential visual contexts in the entire image v . π can either be a deterministic sampling function, or a stochastic sampling process with a probabilistic distribution over Ω . Suppose we acquire n samples v_1, v_2, \dots, v_n according to $\pi(\cdot|v_d)$, we denote the minimum deviation of the resulted token probability from that of the optimal visual context v^* as

$$h_\pi(v^*, n) = \min_{i=1, \dots, n} g(v^*, v_i) = \min_{i=1, \dots, n} D(p_\theta(\cdot|v^*), p_\theta(\cdot|v_i)) \quad (8)$$

where D is the aforementioned symmetric discrepancy measure between two probability distributions, which is within the range of $[0, 1]$. Having a small value of $h_\pi(v^*, n)$ would indicate that we can find some visual context that is close to the optimal v^* through n samples.

We proceed to estimate the minimum deviation $h_\pi(v^*, n)$ from the optimal visual context v^* with n samples. We introduce a partition based on the occurrence of two probabilistic events: the event A where at least one of the samples falls into the ϵ -neighborhood $\mathcal{B}(v^*, \epsilon)$ close to v^* , and its complement. Let us denote the probability of at least one sample falling within $\mathcal{B}(v^*, \epsilon)$ as $P(A)$, and the complementary event's probability as $P(\neg A) = 1 - P(A)$. Hence, we can express the minimum divergence $h_\pi(v^*, n)$ as a marginalization over these events:

$$h_\pi(v^*, n) = P(A) \cdot [h_\pi(v^*, n)|A] + P(\neg A) \cdot [h_\pi(v^*, n)|\neg A] \quad (9)$$

Recognizing that for the one sample in the vicinity of v^* in the event of A , its decoding token probability deviation from the optimal is bounded by $\delta \ll 1$ based on our assumption. Hence we have

$$h_\pi(v^*, n) \leq P(A) \cdot \delta + P(\neg A) \cdot 1 \leq \delta + P(\neg A) \quad (10)$$

Next, we consider two instances of the sampling function $\pi(\cdot|v_d)$ that yield an upper bound for $h_\pi(v^*, n)$.

Normal Distribution Sampling. Suppose sampling from π follows a stochastic process following a normal distribution around v_d . We denote this sampling process as $\pi_g(\cdot|v_d) \sim \mathcal{N}(v_d, \sigma^2 I)$, where we assume a variance of σ^2 for each element of the visual context representation (width, height, center) independently. For $\tilde{v} \in \Omega$, the probability of sampling \tilde{v} following the multivariate normal distribution is

$$q(\tilde{v}; v_d, \sigma^2 I) = \frac{1}{\sqrt{(2\pi\sigma^2)^s}} \exp\left(-\frac{1}{2\sigma^2}(\tilde{v} - v_d)^\top (\tilde{v} - v_d)\right)$$

where $s = 3$ is the dimension of the FOV representation vector. The probability of event $\neg A$ happening, which is none of n FOV samples falling within the ϵ -neighborhood of v^* , is

$$P(\neg A) = P(\|v_1 - v^*\| > \epsilon) \wedge P(\|v_2 - v^*\| > \epsilon) \wedge \dots \wedge P(\|v_n - v^*\| > \epsilon) \quad (11)$$

$$= P(\|\tilde{v} - v^*\| > \epsilon)^n \quad (12)$$

$$= P(\|\tilde{v} - (v_d - \eta)\| > \epsilon)^n \quad (13)$$

From the normal distribution assumption of \tilde{v} , we know that $\tilde{v} - (v_d - \eta)$ also follows a normal distribution $\mathcal{N}(\eta, \sigma^2 I)$. Therefore,

$$P(\neg A) = (1 - P(\|\tilde{v} - (v_d - \eta)\| \leq \epsilon))^n \quad (14)$$

$$= \left(1 - \int_{\nu: \|\nu\| \leq \epsilon} \frac{1}{\sqrt{(2\pi\sigma^2)^s}} \exp\left(-\frac{1}{2\sigma^2}(\nu - \eta)^\top(\nu - \eta)\right) d^s \nu\right)^n \quad (15)$$

$$= (1 - C_g(\epsilon, \eta; \sigma))^n \quad (16)$$

where we use $C_g(\epsilon, \eta; \sigma) \in (0, 1)$ to denote the constant value given ϵ , η , and σ . Following Eq. (10), we now have

$$h_{\pi_g}(v^*, n) \leq \delta + (1 - C_g(\epsilon, \eta; \sigma))^n \quad (17)$$

where the second term goes to 0 as n is increasing to larger values.

Exponential Expansion Sampling. Now suppose sampling from π follows an exponential expanding process, where a sample can be expressed as $v_r = (w_r, h_r, p_r) = ((1 + \lambda)^r w_d, (1 + \lambda)^r h_d, p_d)$ with an expanding factor λ (assuming $\lambda > 0$ without loss of generality) and some r .² Essentially, the sample space comprises all fields of view (FOVs) that maintain the same aspect ratio (i.e. w_d/h_d) and the same center p_d with v_d . Assume the sampling is uniform among all possible FOVs in the sample space, which we denote as $\pi_e(\cdot|v_d) \sim \mathcal{U}(r \in [r_{\min}, r_{\max}])$, where r_{\min} and r_{\max} correspond to the smallest FOV allowed (such as a few pixels) and the largest FOV possible (i.e. the entire original image v), respectively.

For this sampling distribution, we introduce two moderate assumptions regarding the initial detection v_d . First, the center of the detection is relatively close to the optimum, such that $|p_d - p^*| < \epsilon$. Second, The detection v_d and the optimum v^* share the same aspect ratio, meaning $w_d/h_d = w^*/h^*$. This assumption is reasonable since the optimum is unknown, and we can assume it adheres to the aspect ratio used by a standard detector.

We begin by deriving the range of r such that v_r falls into the small neighborhood $\mathcal{B}(v^*, \epsilon)$ around v^* . We need

$$\|v_r - v^*\| \leq \epsilon \quad (18)$$

$$\implies (w_r - w^*)^2 + (h_r - h^*)^2 + (p_r - p^*)^2 \leq \epsilon^2 \quad (19)$$

$$\implies [(1 + \lambda)^r w_d - w^*]^2 + [(1 + \lambda)^r h_d - h^*]^2 + (p_d - p^*)^2 \leq \epsilon^2 \quad (20)$$

⋮

$$\implies (w_d^2 + h_d^2) \left((1 + \lambda)^r - \frac{w_d w^* + h_d h^*}{(w_d^2 + h_d^2)} \right)^2 \leq \epsilon^2 - (p_d - p^*)^2 - \frac{h_d^2 h^{*2}}{(w_d^2 + h_d^2)} \left(\frac{w_d}{h_d} - \frac{w^*}{h^*} \right)^2 \quad (21)$$

$$= \epsilon^2 - (p_d - p^*)^2 > 0 \quad (22)$$

Denoting constants $C_a = \frac{\epsilon^2 - (p_d - p^*)^2}{(w_d^2 + h_d^2)}$ and $C_b = \frac{w_d w^* + h_d h^*}{(w_d^2 + h_d^2)}$, we get the range of r such that $v_r \in \mathcal{B}(v^*, \epsilon)$ as

$$\max\left(r_{\min}, \frac{\log(C_b - \sqrt{C_a})}{\log(1 + \lambda)}\right) \leq r \leq \min\left(r_{\max}, \frac{\log(C_b + \sqrt{C_a})}{\log(1 + \lambda)}\right) \quad \text{if } C_b > \sqrt{C_a} \quad (23)$$

$$\text{Or} \quad r_{\min} \leq r \leq \min\left(r_{\max}, \frac{\log(C_b + \sqrt{C_a})}{\log(1 + \lambda)}\right) \quad \text{if } C_b \leq \sqrt{C_a} \quad (24)$$

We further denote this range as $r \in [C_{\min}(\epsilon, v^*, v_d; \lambda), C_{\max}(\epsilon, v^*, v_d; \lambda)]$, with $r_{\min} \leq C_{\min}(\epsilon, v^*, v_d; \lambda) < C_{\max}(\epsilon, v^*, v_d; \lambda) \leq r_{\max}$. Based on the independent uniform sampling

²Besides expansion, this could also be an exponential shrinking process when r is negative. We abuse the use of ‘‘expansion’’ for both.

assumption, the probability of the event $\neg A$ that none of the n samples fall into the ϵ -neighborhood around the optimum $\mathcal{B}(v^*, \epsilon)$ is

$$P(\neg A) = \left(1 - \frac{C_{\max}(\epsilon, v^*, v_d; \lambda) - C_{\min}(\epsilon, v^*, v_d; \lambda)}{r_{\max} - r_{\min}}\right)^n = (1 - C_e(\epsilon, v^*, v_d; \lambda))^n \quad (25)$$

where we use $C_e(\epsilon, v^*, v_d; \lambda) \in (0, 1]$ to denote the constant value depending on $\epsilon, v^*, v_d, \lambda$. Following Eq. (10), we then have

$$h_{\pi_e}(v^*, n) \leq \delta + (1 - C_e(\epsilon, v^*, v_d; \lambda))^n \quad (26)$$

where the second term goes to 0 as n is increasing to larger values.

Discussion. In the above, we demonstrated that beginning with the initial detected visual context v_d , under certain mild conditions, acquiring n samples according to a distribution $\pi(\cdot|v_d)$ is an efficient method for identifying a sample that leads to a small bounded deviation in the token decoding probabilities from those derived from the optimal visual context v^* . The more samples acquired, the tighter the bound is. This provides a simple and robust way of approximating the optimum.

Different sampling distributions have distinct characteristics. For normal distribution sampling $\pi_g(\cdot|v_d) \sim \mathcal{N}(v_d, \sigma^2 I)$, the variance parameter σ^2 determines the spread of the samples and thus the likelihood of approximating the optimal v^* within $\mathcal{B}(v^*, \epsilon)$. For exponential expansion sampling $\pi_e(\cdot|v_d) \sim \mathcal{U}(r \in [r_{\min}, r_{\max}])$ with samples $v_r = ((1 + \lambda)^r w_d, (1 + \lambda)^r h_d, p_d)$, the parameter λ controls the rate of growth for the sampled visual contexts. In practice, we apply discrete integer values of r to acquire different samples efficiently, thus λ affects the sample coverage of the visual information around v^* .

The choice of the sampling distribution π is contingent upon factors such as the quality of the detector \mathcal{G}_d , the LVLM backbone $\mathcal{M}_\theta^{\text{LVLM}}$, the textual query x , and the visual input v . Specifically, the continuous normal distribution is advantageous for concentrated sampling around v_d , which is particularly effective when the detection perturbation η is small (meaning v_d is near v^*). In contrast, exponential expansion sampling covers an extended range of visual contexts quickly, which is preferable when limited context information is obtained. In scenarios where significant underestimation or overestimation in \mathcal{G}_d detection is present, the exponential expanding strategy can discover the optimal visual context more effectively. \square

E EXPERIMENTS RESULTS

E.1 MORE CHAIR AND POPE RESULTS ON MSCOCO

CHAIR. We also investigate how HALC performs with longer responses, as showed in Fig. 4, where we plot both the number of generated (dashed) and hallucinated (solid) objects with randomly sample 100 images. This experiment is important to further assess HALC’s robustness, as it is commonly believed that OH happens more with objects positioned later in the responses (Zhou et al., 2023), as well as in longer responses (Huang et al., 2023). We observe that HALC is the only method that can keep even smaller number of hallucinations while the number of generated objects increases, demonstrating its superior performance and advantageous robustness in reducing OH.

POPE. As all the numbers in Table 2 are averaged results of the three sampling methods (random, popular and adversarial, as in the original POPE), the complete version of the table is shown in Table 3. HALC outperforms other methods in most of the settings.

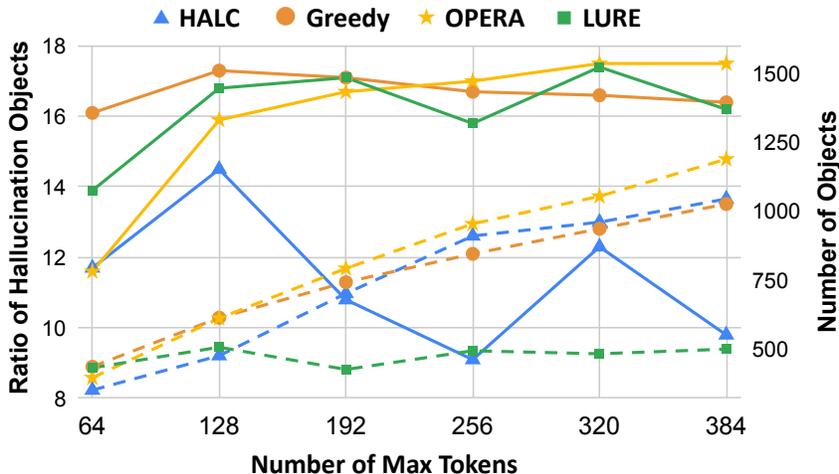


Figure 4: Comparing four mainstream methods on the ratio of hallucination objects ($CHAIR_I$) v.s. the number of max tokens. The right axis (dashed line) indicates the total number of generated objects. HALC outperforms all other methods by maintaining a low ratio of hallucination with the increasing of generated objects.

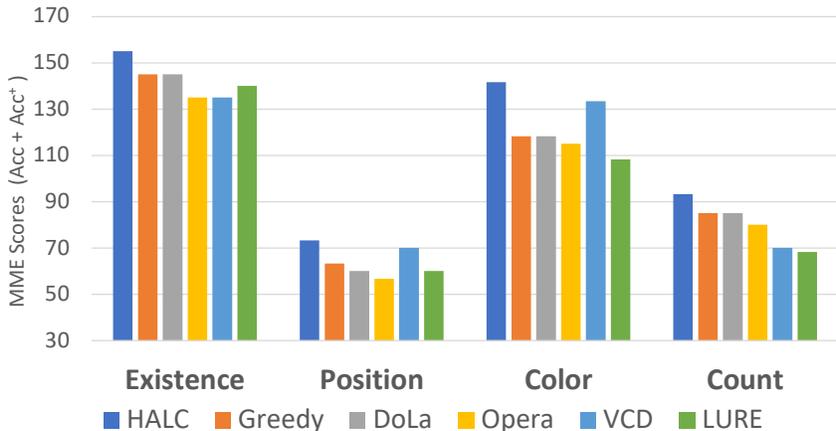


Figure 5: Comparison across OH baselines and SOTAs on four OH-critical MME subsets. All methods adopt MiniGPT-4 as LVM backbone. HALC outperforms all other methods with a large margin: *existence*: +10.7%; *position*: +18.3%; *color*: +19.4% and *count*: +20.2% in average.

E.2 MME

The Multimodal Large Language Model Evaluation (MME) (Fu et al., 2023) benchmark is a comprehensive tool designed to quantitatively compare multimodal LLMs. Following Yin et al. (2023); Leng et al. (2023), we utilize the “existence” and “count” subsets to evaluate the object existence hallucinations and the “position” and “color” subsets for object attribute and relationship hallucination. The comprehensive results across six methods are reported in Fig. 5, where HALC significantly outperforms all the other methods on each sub-task, indicating an overall performance gain in reducing OH while preserving generation quality. The numerical results of MME can be found in Appendix H.

E.3 LLaVA-BENCH QUALITATIVE STUDY

LLaVA-Bench (Liu et al., 2023a) is a collection of 24 images, where each image is paired with a detailed, manually-crafted description and carefully selected questions. The questions are divided into three categories: simple QA (conversation), detailed descriptions, and complex reasoning. In this experiment, we leverage LLaVA-Bench as a case study to qualitatively compare the decoding outputs of HALC with other methods. Results generated by HALC and other OH reduction baselines incorporating mPLUG-Owl2 (Ye et al., 2023), MiniGPT-4 (Zhu et al., 2023; Chen et al., 2023), and LLaVA (Liu et al., 2023b) LVM backbones are shown in Fig. 6, 7 and 8 respectively. In all the plots, red fonts indicate OH, including any of the object existence, attribute or relationship hallucinations.

Table 3: Detailed OPOPE results with random, popular and adversarial samplings.

Setting	Model	Decoding	Accuracy	Precision	Recall	$F_{0.2}$ Score
Random	MiniGPT-4	Greedy	68.30	97.24	37.67	91.67
		Beam Search	68.37	96.30	38.20	90.98
		DoLa	68.50	97.27	38.07	91.78
		OPERA	68.67	96.98	38.53	91.63
		VCD	67.10	96.22	35.60	90.30
		Woodpecker	69.07	96.99	39.366	91.83
		LURE	69.50	96.65	40.4	86.76
		HALC	67.90	97.36	40.4	91.74
	LLaVA-1.5	Greedy	72.20	97.17	45.73	93.14
		Beam Search	71.33	97.48	43.80	93.09
		DoLa	72.30	96.78	46.13	92.86
		OPERA	71.20	96.76	43.87	92.47
		VCD	72.07	96.89	45.60	92.87
		Woodpecker	70.83	95.89	43.53	91.65
		LURE	71.67	97.24	44.6	93.02
		HALC	71.87	97.86	44.73	93.58
	mPLUG-Owl2	Greedy	71.27	96.91	43.93	92.62
		Beam Search	70.50	97.26	42.20	92.61
		DoLa	71.47	96.92	44.33	92.69
		OPERA	70.17	96.92	41.67	92.22
		VCD	70.93	97.31	43.07	92.81
		Woodpecker	70.27	97.99	41.38	93.09
		LURE	70.83	96.71	43.13	92.30
		HALC	71.50	97.38	44.20	93.07
Popular	MiniGPT-4	Greedy	66.43	88.70	37.67	84.30
		Beam Search	67.00	90.09	38.20	85.62
		DoLa	66.8	89.50	38.07	85.08
		OPERA	66.80	88.65	38.53	84.43
		VCD	65.47	65.47	35.60	83.64
		Woodpecker	67.37	89.47	39.37	85.29
		LURE	67.8	89.38	40.4	85.40
		HALC	66.37	90.02	36.80	85.27
	LLaVA-1.5	Greedy	70.27	89.79	45.73	86.58
		Beam Search	69.80	91.25	43.8	87.6
		DoLa	70.43	89.75	46.13	86.60
		OPERA	69.63	90.51	43.87	86.95
		VCD	70.57	91.08	45.60	87.71
		Woodpecker	69.37	90.07	43.53	86.51
		LURE	69.63	89.32	44.6	86.00
		HALC	70.03	90.74	44.67	87.28
	mPLUG-Owl2	Greedy	69.30	89.13	43.93	85.74
		Beam Search	68.83	90.27	42.20	86.48
		DoLa	69.53	89.35	44.33	85.99
		OPERA	69.03	92.02	41.67	87.94
		VCD	69.43	91.10	43.07	87.35
		Woodpecker	68.58	90.73	41.38	86.75
		LURE	69.17	89.99	43.13	86.38
		HALC	69.63	89.95	44.20	86.50
Adversarial	MiniGPT-4	Greedy	65.60	85.35	37.67	81.38
		Beam Search	66.3	87.21	38.20	83.11
		DoLa	65.87	85.74	38.07	81.80
		OPERA	66.3	86.66	38.53	82.68
		VCD	64.77	85.44	35.60	81.08
		Woodpecker	66.88	87.53	39.37	83.60
		LURE	67.13	86.82	40.4	83.14
		HALC	66.00	88.47	36.80	83.94
	LLaVA-1.5	Greedy	69.23	86.30	45.73	83.44
		Beam Search	68.47	86.45	43.8	83.33
		DoLa	69.33	86.07	46.13	83.30
		OPERA	68.37	86.01	43.87	82.95
		VCD	69.37	86.91	45.60	83.99
		Woodpecker	69.20	89.45	43.53	85.96
		LURE	68.7	86.1	44.6	83.13
		HALC	69.87	90.21	44.67	86.80
	mPLUG-Owl2	Greedy	68.73	87.16	43.93	83.98
		Beam Search	68.27	88.17	42.20	84.63
		DoLa	68.87	87.02	44.33	83.91
		OPERA	68.57	90.22	41.67	86.35
		VCD	69.07	89.69	43.07	86.10
		Woodpecker	67.85	87.94	41.38	84.29
		LURE	67.73	84.91	43.13	81.86
		HALC	69.23	88.50	44.20	85.21



HALC: The image is a painting of a face, resembling the famous Mona Lisa's portrait. The painting is done in a modern style, with a blend of colors and textures. The woman's face is the main focus of the painting, and her expression appears to be calm and serene. The background is slightly blurred, adding depth to the artwork.

Greedy: The image is a painting of a woman, possibly a Mona Lisa-style portrait, with a smile on her face. **She is wearing a necklace and appears to be looking directly at the viewer.** The painting is set against a backdrop of a mountain landscape, adding a sense of depth and beauty to the composition.

OPERA: The image features a painting of the famous Mona Lisa, a portrait of a woman attributed to the Renaissance artist Leonardo da Vinci. The Mona Lisa is known for her enigmatic smile and captivating gaze. **The painting is displayed on a wall, showcasing the artwork's beauty and historical significance.**

VCD: The image is a painting of the famous Mona Lisa, capturing the enigmatic smile of the subject. **The artwork is displayed on a wall, and it features a blue background that accentuates the painting's colors and the subject's eyes, making the painting stand out.** The Mona Lisa is the central focus of the artwork, showcasing her famous smile and captivating gaze.



HALC: The image showcases a beautiful view of a mountain range with a large, sandy beach in the foreground. The beach is bustling with activity, with numerous cars and trucks parked along the shore. In the background, the city can be seen, stretching across the landscape. The sky above the scene is filled with clouds, adding to the picturesque atmosphere. The combination of the mountain range, the beach, and the city creates a captivating and dynamic scene.

Greedy: The image showcases a beautiful view of a mountain range with a large, sandy beach in the foreground. **The beach is bustling with activity, with numerous people enjoying the sunny day. Some of the people are standing closer to the water, while others are scattered along the beach, engaging in various activities.** In the background, the mountain range stretches across the horizon, creating a picturesque scene. The combination of the sandy beach, the lush green mountains, and the blue ocean creates a serene and inviting atmosphere.

OPERA: The image showcases a breathtaking view of a mountainous area overlooked by water. A sandy beach stretches along the coast, surrounded on three sides by mountains. **In addition to the beach, there are several boats visible in the water, adding to the picturesque scene. A few people can be seen scattered around the beach, enjoying the beautiful surroundings. Some of them are closer to the water, possibly swimming or sunbathing, while others are further back on the beach.**

VCD: The image depicts a picturesque view of a mountainous landscape with a beautiful blue ocean in the background. The ocean water is visible near the bottom of the image, with a sandy beach extending along the shoreline. The water appears to be calm and inviting, creating a serene atmosphere. The mountainous area is characterized by a series of hills and mountains, with some of them rising steeply and others gently. **A few people can be seen scattered throughout the scene, likely enjoying the natural beauty of the area.**

Figure 6: LLaVA-Bench results comparing HALC and other methods with mPLUG-Ow12 (Ye et al., 2023) backbone.



HALC: The image shows a dining table with a variety of items on it, including a bowl of noodles, a cup of green tea, and some other dishes. The bowl is placed in the center of the table, surrounded by several bowls containing different foods. A cup of green tea is near the top right corner of the table. A spoon can be seen resting on the table, close to the bowl of noodles.

Greedy: The image shows a bowl of noodles with chopsticks and a spoon on the side. The bowl of noodles is filled with a variety of ingredients, including meat, vegetables, and sauce. **The background of the image is a wooden table with a red and white checkered tablecloth. There are also two glasses of beer on the table.**

OPERA: The image shows a bowl of ramen noodles with chicken and vegetables on the side. **The bowl of ramen noodles has a clear broth with chicken and vegetables on the side.** The noodles are cooked and have a chewy texture. **The chicken is cooked and has a crispy texture.** The vegetables are cooked and have a crunchy texture. **The table is made of wood and has a red and white checkered tablecloth on it.**

VCD: The image is a table with a bowl of noodles and some vegetables in it. There are also some chopsticks on the table. **There are two chairs on either side of the table, and a vase of flowers on the table.** The background of the image is a wall with a large window with a view of a city street.



HALC: The image is a toy of a dinosaur with glasses on its face. The dinosaur is made of plastic and has a red and orange flame on its back. The flame is made of plastic and has a red and orange color. The dinosaur is standing on its feet and has a pair of glasses on its face. The dinosaur has a smiling expression on his face.

Greedy: The image is a toy of a dinosaur with glasses on its face. The dinosaur is made of plastic and has a red and orange color scheme. The glasses are made of plastic and have a black frame and lenses. **The dinosaur is standing on its hind legs and has a long tail.** The image is a toy of a dinosaur with glasses on its face. The dinosaur is made of plastic and has a red and orange color scheme.

OPERA: The image is of a toy dinosaur with glasses on its face. The dinosaur is made of plastic and has a red and orange flame on its back. **The dinosaur is standing on its hind legs and has a smile on its face.** The image is of a toy dinosaur with glasses on its face. The dinosaur is made of plastic and has a red and orange flame on its back. The dinosaur is standing on its hind legs and has a smile on its face.

VCD: The image is a toy **flamingo** with a red body and orange wings. **It has a pair of glasses on its head and is standing on its hind legs.** The toy is made of plastic and has a smooth surface. The image is a toy **flamingo** with a red body and orange wings. It has a pair of glasses on its head and is **standing on its hind legs.** The toy is made of plastic and has a smooth surface. The flamingo has a cartoonish appearance, with large eyes and a smiling expression.

Figure 7: LLaVA-Bench results comparing HALC and other methods with MiniGPT-4 (Zhu et al., 2023; Chen et al., 2023) backbone.



HALC: The image features a man standing on the back of a yellow taxi cab, holding a blue white towel. He appears to be ironing clothes while riding in the taxi. The taxi is driving down a busy with other such as a and a car visible in the scene. The man manages to balance himself on the back of the taxi, making the scene quite unique and interesting.

Greedy: The image features a man standing on the back of a yellow taxi cab, holding a blue shirt. He appears to be ironing the shirt while the taxi is in motion. **The taxi is driving down a busy street, with other vehicles such as a truck and a car visible in the scene. Additionally, a traffic light can be seen in the background, indicating that the taxi is driving through an intersection.**

OPERA: In the image, a man is standing on the back of a yellow taxi cab, ironing clothes. He is wearing a yellow shirt and appears to be focused on his task. The taxi is surrounded by other vehicles, including a truck on the left side and a car on the right side of the scene. Additionally, there is a handbag placed near the center of the scene, possibly belonging to the man ironing clothes or someone else in the vicinity.

VCD: The image depicts a busy street scene where a man is **standing on top of a car**, ironing clothes. He is positioned towards the center of the scene, and **there are multiple other people visible in the area.** Various vehicles are present on the street, **including a truck on the left side, a bus on the right side. There are also a few more cars on the right side of the image, near the background.**



HALC: The image features a young boy wearing a suit and tie, with a red bow tie. He is wearing glasses and has a smiling expression on his face. The boy appears to be the main focus of the scene, with the rest of the image being blurry. The suit and tie give him a formal and polished appearance.

Greedy: The image features a young man wearing a suit and tie, with a red bow tie. He is wearing glasses and has a smiling expression on his face. The man appears to be the main focus of the scene, with the rest of the image being blurry. The suit and tie combination, along with the red bow tie, gives the impression that the man is dressed up for a special occasion or event.

OPERA: The image features a young boy wearing glasses and a red bow tie, **standing in front of a building.** He appears to be the main focus of the scene. The boy is wearing a suit and tie, giving him a formal and polished appearance. **The building in the background adds to the overall setting of the image.**

VCD: The image depicts a young boy wearing a blue suit with a red bow tie. He is dressed in a business suit, and his face is adorned with a pair of glasses. The boy has a charming smile and appears to be posing for a picture. **In the background, there is a person partially visible, possibly taking the photo.**

Figure 8: LLaVA-Bench results comparing HALC and other methods with LLaVA (Liu et al., 2023b) backbone.

F DETAILED EXPERIMENTAL SETUPS

The overall experiment settings is reported in Table 4. While the regular greedy decoding follows this setting, the beam search variant in our experiment essentially applies a token-wise beam search based on accumulated probability scores of the previous tokens $y_{<t}$. We use the default code for implementation of these two baselines in HuggingFace TransformersRepository (Wolf et al., 2020).³

Table 4: Overall Experiment Settings

Parameters	Value
Maximum New Tokens (CHAIR)	64
Maximum New Tokens (POPE)	64
Maximum New Tokens (MME)	128
Top-k	False
Top-p	1
Temperature τ	1

The complete hyper-parameters for HALC in our experiments in §4 is reported in Table 5. Specifically, there are four major hyper-parameters that can actively adjust the effectiveness of HALC to adapt to different task settings:

1. *FOV Sampling Distribution*: Typically, a normal distribution, which concentrated around v_d , provides a tighter bound under minimal perturbations, while an exponential distribution, with a more averaged coverage of the sampling space, is preferable when less contexts of the task is available. Thus to preserve generality in our experiment, we have employed the exponential distribution with exponential growth factor $\lambda = 0.6$.
2. *Number of Sampled FOVs n* : n determines the number of sampled FOVs in the discretized FOV space. According to Theorem D.1, while increasing n and adjusting the distribution parameters can efficiently reduce C_S and enhance the robustness against bounded perturbations, it’s notable that the runtime costs also raise with n . Consequently, we set $n = 4$ across all our experiments.
3. *JSD Buffer Size m* : For each beam in the overall beam search process (beam size k), our bi-adaptive visual grounding module samples n visual contexts, which through interpolated JSD calculation would produce $\frac{n \cdot (n-1)}{2}$ JSD values in total. Then we select the top m FOV pairs with relatively large discrepancy to produce contrastive candidate distributions.
4. *Beam Size k* : The beam size k is set to adjust the diversity and range for HALC to search for the best candidate captions. Essentially, the global visual matching score module selects the top k diverse captions from $2m \cdot k$ text sequence candidates passed from the local adaptive visual grounding module. While a larger k involves a larger search space and hopefully a better generation, the runtime cost also raises linearly w.r.t. k . HALC adopts Bootstrapping Language-Image Pre-training (BLIP) (Li et al., 2022a) for both text and image encoding when computing their cosine similarity scores. Notably given the global search capability of our visual matching score module, HALC seeks to preserve a more diverse set of captions within the beam buffer.

Table 5: HALC Hyperparameter Settings

Parameters	Value
Amplification Factor α	0.05
JSD Buffer Size m	6
Beam Size	1
FOV Sampling	Exponential Expanding
Number of Sampled FOVs n	4
Exponential Growth Factor λ	0.6
Adaptive Plausibility Threshold	0.1

³<https://huggingface.co/docs/transformers>

Regarding the comparison of HALC with SOTAs that are specifically designed for OH mitigation, we adopt the code, hyper-parameters, and pre-trained models of each method outlined in their public repositories and papers respectively. Specifically, the hyper-parameters for DoLa (Chuang et al., 2023)⁴ is reported in Table 6; OPERA (Huang et al., 2023)⁵ is reported in Table 7; and the hyperparameters for VCD (Leng et al., 2023)⁶ is reported in Table 8. For each of these baselines, we strictly follow their implementations and hyper-parameters as reported in the paper to reproduce their results.

Table 6: DoLa Hyperparameter Settings

Parameters	Value
Repetition Penalty θ	1.2
Adaptive Plausibility Threshold β	0.1
Pre-mature Layers	[0, 2 ··· , 32]

Table 7: OPERA Hyperparameter Settings

Parameters	Value
Self-attention Weights Scale Factor θ	50
Attending Retrospection Threshold	15
Beam Size	3
Penalty Weights	1

Table 8: VCD Hyperparameter Settings

Parameters	Value
Amplification Factor α	1
Adaptive Plausibility Threshold	0.1
Diffusion Noise Step	500

Regarding post-hoc correction method woodpecker (Yin et al., 2023)⁷ and LURE (Zhou et al., 2023)⁸, we also strictly follow their implementations and hyper-parameters settings as reported in the corresponding papers to reproduce their results. For woodpecker, we adopt their original code and use OpenAI API to access GPT-3.5 Turbo. In average, per 500 images would result in approximately \$4.5 cost. For LURE, we also directly adopt their pre-trained projection layer model (for MiniGpt4) to reproduce the results reported in this paper. All the hyper-parameters are default.

G EMPIRICAL STUDIES ON OPTIMAL VISUAL CONTEXTS

We verify our insight that optimal visual context is important in correcting object hallucination through an empirical pilot study. Fig. 1 shows the oracle performance of OH levels when we rely on optimal visual contexts for tokens through brute-force search, with greedy decoding on the MME benchmark (Fu et al., 2023) on three categories of OH. Specifically, each MME sub-task contains 30 images, and we have followed (Leng et al., 2023) and selected four sub-tasks (including *existence*, *count*, *color*, *position*) to evaluate the hallucination in our analysis, in total 110 distinct images. Based on these images, we manually constructed multiple challenging questions (2-4 per image) that are likely to induce the LVM to hallucinate (e.g. some minor objects in the distance or some plausible but unfaithful objects that are likely to co-occur). Then we take each question as a count unit and calculate the number of hallucinations on word level (instead of token level) which could attributed for each of the three main OH sources. Then for each question with a hallucination occurring, we first search across the original image input using a brutal-force breadth-first algorithms until the hallucinating token is corrected to be consistent with the ground truth. This process effectively

⁴<https://github.com/voidism/DoLa>

⁵<https://github.com/shikiw/OPERA>

⁶<https://github.com/DAMO-NLP-SG/VCD>

⁷<https://github.com/BradyFU/Woodpecker>

⁸<https://github.com/YiyangZhou/LURE>

succeed to retrieve the optimal visual context for 54.0% of the questions. For those questions that fail this brutal-force search, we further manually select the visual context candidates based on human priors. In total, 84.5% of the questions that contain these three sources of hallucinations can be eliminated with an explicit optimal visual prior v^* .

H MME EXPERIMENT DETAILS

The experiment details mostly follow Appendix G, where we adopt each sub-task of 30 images from MME dataset, and reconstruct the question prompt following OPOPE. Specifically, instead of simply asking a yes/no question, we first ask the decoder to generate a detailed caption then check whether the MME-targeted positive/negative word existed in the caption. While this modified offline metric could lower the score for false negative, we argue this holds fair across all methods. The detailed results are reported in Table 9.

Table 9: Comparison of Decoder Performances on 4 MME sub-tasks

Decoder	Existence	Position	Color	Count	Max Tokens	Num of Samples
HALC	155	73.33	141.67	93.33	128	110
Greedy	145	63.33	118.33	85	128	110
DoLa	145	60	118.33	85	128	110
Opera	135	56.67	115	80	128	110
VCD	135	70	133.33	70	128	110
LURE	140	60	108.33	68.33	128	110