# HALLUTEXT: TOWARDS BENCHMARKING AND MITI-GATING OCR HALLUCINATION FOR LVLMS

**Anonymous authors**Paper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

033

035

037

040

041

042

043

044

046

047

048

051

052

# **ABSTRACT**

Optical Character Recognition (OCR) serves as a critical bridge connecting vision and language, attracting increasing attention in the community of Large Vision-Language Models (LVLMs). However, due to the prevalent encode-then-decode architecture, LVLMs tend to over-rely on language priors, leading to frequent failures in following basic visual-text instructions. We term this issue OCR hallucination. To systematically mitigate it and facilitate reliable OCR perception in LVLMs, we conduct the first large-scale empirical analysis based on OCRBench v2. Our findings reveal that current LVLMs frequently misinterpret or ignore textual visual content, particularly across two orthogonal dimensions, including perception task and hallucination taxonomy. Building on these insights, we introduce HalluText, a benchmark specifically designed to comprehensively evaluate OCR hallucination in LVLMs across nine subclasses. Alongside this benchmark, we propose OCRAssistor, a lightweight plug-and-play method pioneering largesmall model collaboration. By integrating compact OCR model outputs into the LVLM decoding process, it achieves a 9.6% improvement on HalluText with only marginal computational cost. When applied to OCRBench v2, this method also improves the performance of the top-performing open-source model Qwen2.5-VL-7B, achieving a 3% gain and highlighting the importance of addressing OCR hallucination in LVLMs. Through our benchmark and proposed solution, we hope to shed light on the challenges and potential pathways for improving visual text perception in LVLMs. The organized benchmark and the relevant code will be released soon.

# 1 Introduction

Driven by advances from both academia and industry, Large Vision Language Models (LVLMs) are increasingly applied across a wide range of domains. As a crucial bridge between vision and language, Optical Character Recognition (OCR) has emerged as both a foundational pre-training paradigm and a key task for supervised fine-tuning. OCR-centric tasks have also garnered significant attention from both general-purpose (Wang et al., 2024b; Li et al., 2024a; Lu et al., 2024; Yao et al., 2024; Bai et al., 2025; Zhu et al., 2025a) and OCR-specialized LVLMs (Li et al., 2024b; Huang et al., 2024a; Yu et al., 2024b; Zhao et al., 2024; Nacson et al., 2025; Li et al., 2025), owing to their wide applicability in real-world scenarios such as smart offices, content moderation, and document intelligence.

Despite this growing focus on OCR-centric tasks, we observe that current LVLMs still often struggle with seemingly simple questions that involve understanding text within images. Figure 1 illustrates three representative failure cases where state-of-the-art models (Yao et al., 2024; Huang et al., 2024a; Bai et al., 2025) consistently fail. Borrowing the concept of hallucinations, we attribute this issue as "OCR hallucination", defined as instances where the responses generated by LVLMs fail to accurately follow visual text-centered instructions. To systematically analyze and attribute these errors, we conduct a comprehensive empirical study on widely adopted OCRBench v2 (Fu et al., 2024). Evaluations across over 1,000 samples reveal that these errors are prevalent across different LVLMs and tasks, while exhibiting consistent patterns that enable their categorical grouping. Driven by such findings, we categorize the errors along two orthogonal dimensions: 1) the perception task stage, which focuses on the perceptual stages of localization and recognition, and 2) the hallucination taxonomy, which classifies error types into category, relation, and attribute hallucinations. Take the

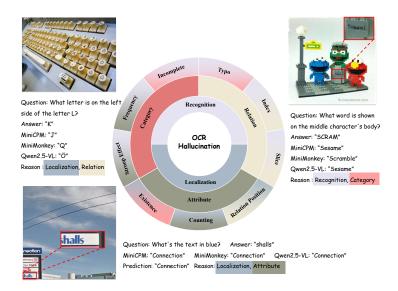


Figure 1: Taxonomy of OCR Hallucinations. The inner rings represent a dual-perspective taxonomy by task stage and hallucination type, while the outer ring indicates the nine HalluText subsets aligned with these categories. Colors denote hallucination types, and three examples around the perimeter illustrate their definitions and characteristics.

the bottom of Figure 1 as an example, the question is "What's the text in blue?", the correct answer is "shalls", but all three models erroneously output "Connection". We attribute this error primarily to incorrect localization—the models fail to attend to the region containing the blue text. Further analysis reveals that this localization failure stems from a misunderstanding of the visual attribute "blue", indicating an error due to attribute hallucination. Thus, this case exemplifies how an initial perceptual failure (e.g., color misinterpretation) can propagate to semantic-level hallucinations. By framing OCR errors within this dual-perspective framework, we aim to better understand the underlying causes of OCR failures and provide actionable insights for future model development.

Based on these insights, we introduce a new benchmark **HalluText**. Unlike the scattered hallucination-related samples in OCRBench v2, HalluText offers a more comprehensive and structured diagnosis for OCR hallucination in LVLMs. Following (Yin et al., 2024), we formulate test samples as multiple-choice questions and collect 4,678 image–question–answer triplets covering nine distinct types of hallucinations. As illustrated in Figure 1, these types are built upon the two previous orthogonal dimensions, comprising existence, incompletion, typo, position, index, slice, counting, frequency, and stroop effect.

Furthermore, we also propose a lightweight, plug-and-play method to mitigate OCR hallucinations, called **OCRAssistor**. This framework pioneers a novel collaborative paradigm between large and small models, where a small-scale OCR-specialized model injects vision-grounded cues to guide the decoding process of large vision-language models (LVLMs). Notably, this design does not require additional fine-tuning of the large model, making it both efficient and flexible to deploy. Despite its simplicity, OCRAssistor achieves impressive results, improving the baseline Qwen2.5-VL-7B by 9.6% on HalluText. When applied to the more general OCRBench v2, it outperforms the baseline by 2.5% on the English subset and 3.7% on the Chinese subset. These results not only demonstrate the effectiveness and generalizability of our approach in fine-grained perception tasks but also underscore the critical importance of addressing hallucination in OCR-centric applications. Our contributions are summarized into three main aspects.

- We conduct an extensive empirical study to uncover the overlooked problem of OCR hallucination in LVLMs. Driven by the results, we establish a dual-perspective taxonomy based on the task categories and hallucination types to systematically analyze these errors.
- We construct HalluText, a fine-grained benchmark for OCR hallucination. HalluText consists of 4,678 carefully curated samples across 9 subsets, each targeting specific perception and hallucina-

tion dimensions. Compared to existing benchmarks like OCRBench v2, HalluText offers a more comprehensive and structured diagnostic of OCR hallucination, providing clear insights for future advancements of LVLMs.

We design OCRAssistor, a plug-and-play method to mitigate OCR hallucination through a novel
large-small model collaboration framework. To our knowledge, it is the first work to adopt such
a collaborative paradigm for OCR hallucination mitigation. OCRAssistor incorporates minimal
computational overhead, yet significantly improves LVLMs on both HalluText and OCRBench v2.
Extensive experiments validate its effectiveness, efficiency, and scalability across a wide range of
scenarios.

# 2 RELATED WORKS

#### 2.1 OCR-AWARE BENCHMARK IN LVLM ERA

Before the LVLM era, OCR-aware benchmarks focused on specific sub-tasks, such as scene text detection and recognition (e.g., ICDAR (Karatzas et al., 2013), Total-Text (Ch'ng & Chan, 2017), SCUT-CTW1500 (Liu et al., 2017)), visual text understanding (e.g., TextVQA (Singh et al., 2019), STVQA (Biten et al., 2019)), key information extraction (e.g., FUNSD (Jaume et al., 2019), SROIE (Huang et al., 2019)), and chart understanding (e.g., ChartQA (Masry et al., 2022), infographicVQA (Mathew et al., 2022)). With the rise of LVLMs, the focus shifted towards unified OCR-centric benchmarks. OCRBench (Liu et al., 2024) integrates five major tasks—text recognition, scene VQA, document VQA, key information extraction, and handwritten formula recognition—across 27 datasets. The latest OCRBench v2 expands further, adding element parsing, knowledge reasoning, and mathematical calculations. Additionally, document parsing and understanding have gained widespread attention, with new benchmarks (Wei et al., 2024; Ouyang et al., 2025; Li et al., 2025) created for evaluating document-specific tasks. Recent research has also analyzed segmentation deficiencies, with OCR-Reasoning (Huang et al., 2025) and Reasoning OCR (He et al., 2025a) focusing on dense text understanding. Work by (Shu et al., 2025) and (He et al., 2025b) addresses hallucinations in non-semantic and occluded/blurred text-rich scenarios. In this paper, we propose a new benchmark to uncover OCR hallucinations in OCR-centric tasks, based on common failure cases from the general OCRBench v2.

# 2.2 HALLUCINATION MITIGATION

The concept of hallucination originates from the domains of pathology and psychology, where it is defined as the perception of something that does not exist in reality (Macpherson & Platchias, 2013). In natural language processing, hallucination typically refers to instances where generated content is implausible or inconsistent with the source input (Maynez et al., 2020). In the LVLM scenario, Hallucination refers to the phenomenon where the generated text response does not align with the corre-

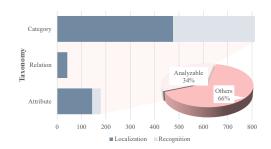


Figure 2: The distribution of failure cases on OCRBench v2.

sponding visual content (Bai et al., 2024). To address this, some methods improve training data (Liu et al., 2023; Yu et al., 2024a; Zhang et al., 2024), adopt architectural designs (Li et al., 2024b; Jain et al., 2024), or introduce post-training stages (Sun et al., 2023; Zhao et al., 2023; Gunjal et al., 2024). Given the high computational cost, others have explored training-free approaches (Leng et al., 2024; Wang et al., 2024c;a; Huang et al., 2024b; Favero et al., 2024; Zhu et al., 2025b), mainly categorized as contrastive decoding and attention intervention. Contrastive decoding (Leng et al., 2024; Wang et al., 2024c;a; Ghosh et al., 2025) modifies the decoding distribution but requires additional inference, leading to latency. Attention intervention (Huang et al., 2024b; Zhu et al., 2025b; Favero et al., 2024)shifts focus toward visual inputs during decoding but still incurs overhead. In OCR-centric tasks, we propose a large-small model collaboration framework that in-

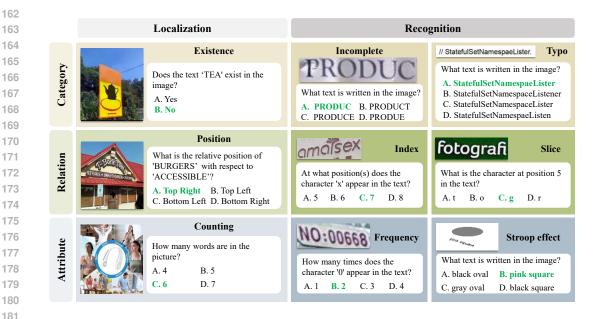


Figure 3: An overview of HalluText, which collects the challenging issues in visual text perception, including localization and recognition. The column axis involves the three categories of hallucination, *Category Hallucination*, *Relation Hallucination*, and *Attribute Hallucination*, respectively. Bold indicates the correct answer.

tegrates a lightweight, visually faithful OCR model into LVLMs to mitigate hallucination. This plug-and-play design improves visual alignment while avoiding the latency of prior training-free methods.

# 3 ANALYSIS ON VISUAL TEXT HALLUCINATION

# 3.1 EMPIRICAL ANALYSIS ON OCRBENCH V2

OCRBench v2 is currently the largest and most comprehensive benchmark for OCR-related tasks, comprising over 10,000 annotated question-answer pairs across more than 20 diverse scenarios in both Chinese and English. However, hallucinated samples are scattered across various task categories, making it challenging to systematically evaluate the hallucination-handling capabilities of LVLMs. To address this, we perform targeted analysis by identifying and categorizing hallucination types in model failures.

We select three representative LVLMs, including Qwen2.5-VL, MiniCPM-o 2.6, and MiniMonkey, and apply the official evaluation scripts to identify 3,006 samples where all models fail consistently. It is important to note that not all error cases are related to hallucinations. We consider only hallucination-related samples as analyzable. Other errors, such as those caused by question misinterpretation, complex reasoning failures, or annotation issues within the dataset, are not directly related to OCR hallucinations and are therefore excluded from our analysis.

To attribute hallucinations, we adopt two orthogonal dimensions: (1) perception task type (localization vs. recognition) and (2) hallucination type, which we refine for OCR settings as follows: Category Hallucination: incorrect recognition of text content or coordinates; Relational Hallucination: errors in spatial or semantic relations between text instances; Attribute Hallucination: incorrect description of text attributes such as quantity or color. Among the failed samples, 1,034 (34%) are deemed analyzable. We perform detailed attribution across models, and the resulting hallucination distributions, which are summarized in Figure 2, reveal consistent patterns across task types and models. These findings provide both theoretical grounding and empirical basis for developing robust hallucination benchmarks in OCR-focused LVLM evaluation.

Table 1: Distribution and original source of HalluText.

Subsets	Existence	Position	Counting	Stroop	Туро	Incomplete	Freq.	Index	Slice
Number	250	500	233	200	500	1495	500	500	500
Source	SCUT-ENS	Total-Text	ICDAR2013	Manual	Typo-corpus	Union-Incomplete	Unio	n-contex	tless
Source	(Liu et al., 2020)	(Ch'ng & Chan, 2017)	(Karatzas et al., 2013)	-	(Hagiwara & Mita, 2019)	(Jiang	et al., 20	23)	

# 3.2 HALLUTEXT BENCHMARK

Building on the empirical analysis in the previous section, we identify the distribution of hallucinated samples within OCRBench v2. Based on the occurrence scenarios of these hallucinations, we construct a dedicated dataset for OCR hallucination research, named HalluText, by reorganizing existing OCR datasets according to hallucination types. HalluText consists of 9 subsets, each corresponding to a specific hallucination category, and includes a total of 4,678 image—question—answer triplets. The definitions and construction procedures of each subset are detailed in the following section. The distribution and sources of the subsets are summarized in Table 1. The detailed construction procedures of all subsets are provided in Appendix B.

**Existence.** Due to training data biases, LVLMs are prone to hallucinations when presented with manipulated images. This subset is constructed from the scene text erasure dataset SCUT-ENS (Liu et al., 2020), with the goal of evaluating whether LVLMs can accurately perceive the presence of specific words in an image. To address the balance between *Yes* and *No* answers, we also incorporate negative polarity questions during the question construction process.

**Incompletion. & Typo** Influenced by the Language model, LVLMs tend to replace non-semantic words in their outputs with semantically plausible text. We constructed the Incomplete and Typo subsets using scene text that is affected by occlusion or truncation, and text containing common spelling errors, respectively. These subsets are designed to evaluate the ability the capability of accurately recognizing visual text while remaining robust to linguistic priors.

**Position** Empirical studies reveal that LVLMs exhibit limitations in relative position perception. To evaluate their ability to understand spatial relationships in real-world scenes, we design a relative position recognition task based on scene text data. This task assesses how well LVLMs can perceive relative positions of visual elements across the entire image.

**Index & Slice** Correspondingly, we also observe relation-level hallucinations within individual text instances. To minimize the influence of semantic priors, we construct position-specific questions on the Union14M-Contextless subset (Jiang et al., 2023), using common string slicing and indexing operations for naming. These subsets are designed to evaluate the ability to perceive intra-word spatial relations within single text instances.

Counting Empirical results suggest that LVLMs struggle with counting-related tasks. To evaluate their ability to perceive numerical attributes of visual text, we construct a counting task based on the ICDAR2013 dataset, which primarily consists of focused text with minimal ambiguity. This subset is designed to assess whether LVLMs can accurately determine the number of text instances in an image.

**Frequency & Stroop Effect** Beyond counting, color perception represents another key aspect of attribute-level hallucination. In addition to the intra-text counting task, we draw inspiration from the Stroop Effect (MacLeod, 1991) to construct synthetic images containing color and shape words. These two subsets are designed to evaluate the ability of LVLMs to suppress hallucinations related to text quantity and text color, respectively. More details are illustrated in Algorithm 1.

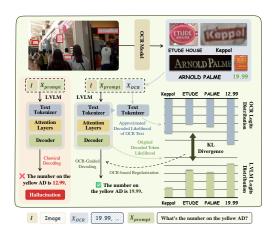
# 3.3 METRIC

Following OCRBench v2, we adopt a multiple-choice QA format with up to four options per question, using accuracy as the evaluation metric. For standardized evaluation, given images  $\mathcal{I}$ , questions  $\mathcal{Q} = \{q_i\}_{i=1}^N$ , and answers  $\mathcal{A} = \{a_i\}_{i=1}^N$ , we employ a fixed prompt template: "Please strictly follow these rules: Only output the letter of the correct answer. Place the answer on a separate last line. Question:  $\{question\}$ . Answer: $\{\}$ ." The templatized questions are then fed into LVLMs to obtain predictions  $\mathcal{P} = \{p_i\}_{i=1}^N = \mathcal{M}(\mathcal{I}, \mathcal{Q})$ , where  $\mathcal{M}$  is the LVLM. The multi-choice accuracy is formulated as  $Acc = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(p_i = a_i)$ , where  $\mathbb{1}$  is the indicator function.

# 4 METHOD

Inspired by (Ghosh et al., 2025), we introduce the OCRAssistor, an OCR-guided decoding framework that pioneers a collaborative mechanism between a large vision-language model (LVLM) and a lightweight OCR expert model to alleviate perception hallucination in OCR-centric tasks. The overall pipeline is illustrated in Figure 4.

Given an input image and a textual prompt, we first extract textual elements from the image using a lightweight OCR model, resulting in a sequence  $X_{\text{OCR}} = \{o_1, o_2, \dots, o_k\}$ . To ground the LVLM's generation in these visual texts, we prepend them to the user prompt  $X_{\text{prompt}} = \{x_1, x_2, \dots, x_n\}$ , yielding the augmented prompt  $X_{\text{concat}} = \{o_1, \dots, o_k, x_1, \dots, x_n\}$ .



Next, we construct a reference token distribution over the model's vocabulary  $\mathcal V$  from the OCR outputs. Specifically, we pass the OCR text through the LVLM's embedding and output layers to obtain pseudo-logits  $\hat\ell_{OCR} \in \mathbb R^{|\mathcal V|}$ , which approximate the token likelihoods if the model were asked to generate the OCR text. These logits are normalized with a temperature-scaled softmax:

$$\hat{p}_{OCR}(w) = \frac{\exp(\hat{\ell}_{OCR}(w)/T)}{\sum_{v \in \mathcal{V}} \exp(\hat{\ell}_{OCR}(v)/T)}, \quad T > 0,$$
(1)

where T controls the sharpness of the distribution. This step ensures that all tokens receive non-zero probability mass. During decoding, let  $\mathcal{L}_i \in \mathbb{R}^{|\mathcal{V}|}$  denote the original decoded logits predicted by the LVLM  $\mathcal{M}$  at step i for the

Figure 4: The framework of OCRAssistor.

next token, and let  $p_i(w) = \operatorname{softmax}(\mathcal{L}_i)(w)$  be the corresponding token distribution. Motivated by KL-divergence (Kullback, 1951), we incorporate the OCR guidance by directly modifying the logits in a distribution-aware manner:

$$\mathcal{L}'_{i}(w) = \mathcal{L}_{i}(w) - \lambda \cdot \log \frac{p_{i}(w)}{\hat{p}_{OCR}(w)}, \quad \forall w \in \mathcal{V},$$
(2)

where  $\lambda$  is a hyperparameter controlling the strength of OCR-based regularization. To ensure the robustness of our design choices, we provide a detailed ablation on  $\lambda$  in Table 10 of Appendix D. Intuitively, tokens more consistent with OCR-derived probabilities are relatively boosted, while inconsistent ones are penalized. The adjusted logits are then normalized to obtain the final decoding distribution:

$$p'(w \mid x_{< i}) = \operatorname{softmax}(\mathcal{L}'_i)(w). \tag{3}$$

At each step, the next token is sampled from p', using the same decoding configuration as the base model. This ensures that the OCR guidance is seamlessly integrated into standard LVLM decoding while encouraging semantic consistency with visual texts. The generation terminates once an end-of-sequence token is produced or a predefined length limit is reached.

This approach integrates the predictions of an external OCR model into the decoding process via a KL-divergence-based guidance mechanism. This alignment encourages the LVLM to focus more heavily on visually grounded textual cues, effectively suppressing hallucinations and improving visual fidelity. Unlike prior comparison-based decoding methods, such as Visual Description Grounding Decoding (VDGD), which require multiple rounds of model inference, our approach achieves efficient inference by performing only one forward pass through the LVLM and a lightweight OCR model. This significantly reduces computational cost while maintaining strong performance.

Table 2: Perfomance comparison on HalluText. OA indicates the OCRAssistor. Abbreviations: EX = Existence, RP = Relative Position, CT = Counting, ST = Stroop Effect, TY = Typo, IC = Incompletion, FQ = Frequency, ID = Index.

Model	Lo	ocalizati	on	4.00			Recog	nition			100	100
Wiodei	EX	POS	CT	$Acc_{loc}$	ST	TY	IC	FQ	ID	SL	$Acc_{rec}$	$Acc_{all}$
	•			Prop	rietary	LVLM	5					
Gemeni-Pro	91.6	63.8	69.1	74.8	99.0	66.8	56.1	63.7	46.2	62.8	65.8	68.8
GPT-40	89.2	49.0	56.7	65.0	98.5	95.2	75.1	64.9	66.9	60.1	76.8	72.8
Open-source LVLMs												
Qwen2.5-VL-3B	90.8	39.4	46.4	58.9	94.0	79.2	42.1	59.6	45.0	47.7	59.6	59.3
Qwen2.5-VL-7B	98.4	50.4	59.7	69.5	90.5	88.8	72.3	62.3	50.3	49.8	69.0	69.1
Qwen2.5-VL-32B	94.4	55.4	63.1	71.0	92.0	88.4	81.6	69.6	44.3	47.3	70.5	70.7
InternVL3-2B	75.2	31.8	44.2	50.4	99.0	82.8	62.6	45.6	38.4	40.7	61.5	57.8
InternVL3-8B	89.2	36.8	67.0	64.3	59.8	90.6	63.4	53.3	52.6	56.6	62.7	63.3
InternVL3-14B	95.6	67.2	63.1	75.3	99.5	92.8	88.1	75.7	60.0	63.8	80.0	78.4
MiniCPM2.6-o-8B	77.6	51.2	51.9	60.2	96.5	82.8	77.1	51.9	44.8	51.4	67.4	65.0
MiniMonkey-2B	76.4	32.4	35.2	48.0	76.9	71.0	77.1	43.8	6.0	24.4	49.9	49.2
LLaVA-NeXT-7B	74.0	38.6	31.3	48.0	71.4	48.6	41.5	64.3	51.8	39.6	44.6	45.7
LLaVA-NeXT-7B + OA	81.6	37.2	42.1	53.6	94.0	63.4	73.5	49.3	32.4	35.7	58.1	56.6 (+10.9)
Qwen2.5-VL-7B	98.4	50.4	59.7	69.5	90.5	88.8	72.3	62.3	50.3	49.8	69.0	69.1
Qwen 2.5-VL-7B+OA	98.8	50.4	66.5	71.9	95.5	89.8	81.3	71.0	72.9	82.0	82.1	78.7 (+9.6)

# 5 EXPERIMENTS

# 5.1 SETTINGS

We evaluate and compare HalluText and OCRBench v2 with several state-of-the-art LVLMs, including proprietary models GPT-40 (Hurst et al., 2024) and Gemini-Pro (Team et al., 2024), as well as open-source models InternVL3 (Zhu et al., 2025a), Qwen2.5-VL (Bai et al., 2025), LLaVA-NeXT (Li et al., 2024a), MiniCPM2.6-o (Yao et al., 2024), and MiniMonkey (Huang et al., 2024a). We facilitate the widely used OCR engine PaddleOCR-v5 as our OCR model. To ensure fair comparison, we locally re-infer representative open-source LVLMs using only the annotated question prompts, and follow the official evaluation protocol. The maximum number of generated tokens is set to 1024. The temperature factor T and regularization factor  $\lambda$  are set to 0.1 and 0.5 by default. The detailed prompt settings are discussed in Appendix C.

# 5.2 RESULTS AND ANALYSIS

## 5.2.1 HALLUTEXT

Table 2 presents the results on our proposed HalluText. We have several findings:

- 1) OCR-centric hallucination remains an unsolved challenge across both proprietary and open-source models. Overall, all models achieve less than 80% accuracy on our benchmark, and their performance on fine-grained subsets, such as relative position, counting, index, and slice, is significantly lower than the average accuracy  $Acc_{all}$ . This highlights the persistent and under-addressed issue of OCR hallucination in current LVLMs. The poor performance on Slice and Index suggests a limited understanding of ordinal relationships. Counting and Relative position tasks remain difficult due to the insensitivity of LVLMs to object-level correlation and the lack of fine-grained perceptual reasoning. Moreover, subsets like Frequency and Slice, which lack contextual information, expose the reliance of LVLMs on semantic cues for accurate recognition.
- 2) The scaling law continues to hold for the OCR hallucination task. We evaluate recent versions of Qwen and InternVL across three model scales and observe a consistent trend: larger models exhibit stronger capabilities in suppressing textual hallucinations, confirming the applicability of scaling effects in this domain.
- 3) Our OCRAssistor method, under a training-free setting, integrates an off-the-shelf open-source OCR model and yields substantial improvements. Specifically, it improves LLaVA by 12% and Qwen2.5-VL by 9.6%, with consistent gains across nearly all fine-grained subsets. These results demonstrate the effectiveness of our approach in mitigating OCR hallucinations.

Table 3: Perfomance comparison on OCRBench v2. OA indicates the OCRAssistor. Abbreviations: TR = Text Recognition, TD = Text Detection, TS = Text Spotting, RE = Relation Extraction, EP=Element Parsing, MC = Metathetical Calculating, TU=Text Understanding, KR = Knowledge Reasoning.

Model				Englis	sh Part				Chinese Part						Overall		
Wiodei	TR	TD	TS	RÉ	EP	MC	TU	KR	TR	RE	EP	TU	KR	English	Chinese		
Proprietary LVLMs																	
Gemini-Pro	61.2	39.5	13.5	79.3	39.2	47.7	75.5	59.3	52.5	47.3	30.9	51.5	33.4	51.9	43.1		
GPT-40	61.2	26.7	0.0	77.5	36.3	43.4	71.1	55.5	21.6	53.0	29.8	38.5	18.2	46.5	32.2		
Open-source LVLMs																	
MiniMonkey-2B	58.1	19.6	0.0	51.3	33.0	15.7	61.7	44.8	61.4	40.5	27.9	42.8	17.9	35.5	38.1		
MiniCPM-o-2.6-8B	67.4	26.5	0.0	70.1	34.0	31.7	70.6	57.6	54.7	52.4	27.6	42.5	31.6	44.8	41.7		
InternVL3-8B	66.9	25.7	0.0	85.3	36.8	34.4	72.3	58.8	67.6	56.9	32.7	53.8	36.7	47.5	49.5		
LLaVA-NeXT-7B	38.0	18.5	0.0	21.0	9.8	13.3	65.9	48.6	5.8	9.3	14.1	4.0	1.6	26.9	7.0		
LLaVA-NeXT+OA	47.2	19.1	0.0	60.4	22.7	22.0	64.4	45.0	31.0	29.1	18.2	44.0	18.1	35.1 (+8.2)	28.1 (+21.1)		
Qwen2.5-VL-7B	67.0	22.3	0.0	76.8	28.2	34.1	72.0	56.3	69.0	52.7	42.3	43.3	37.9	44.6	49.1		
Qwen2.5-VL-7B+OA	60.4	22.6	0.0	86.4	33.6	46.2	72.9	54.7	57.0	64.8	39.4	56.8	45.8	47.1 (+2.5)	52.8 (+3.7)		

# 5.2.2 OCRBENCH V2

We further evaluate our approach on OCRBench v2, a general benchmark for OCR-centric tasks. Table 3 shows that our method improves LLaVA-NeXt-7B by 8.2% in English scenarios and 21.1% in Chinese scenarios. For Qwen2.5-VL-7B, which possesses stronger baseline capabilities, our method still achieves 2.5% and 3.7% improvements in English and Chinese settings, respectively. The notably larger gain in LLaVA's Chinese performance is primarily due to the relatively limited Chinese data exposure during its pretraining phase, compared to Qwen2.5-VL. This suggests that our method can effectively compensate for underrepresented modalities or languages in pretraining, particularly in low-resource scenarios. Beyond perception tasks, OCRAssistor also improves relation extraction, text comprehension, and knowledge reasoning. These results show that integrating an OCR model not only benefits visual-text perception tasks but also enhances high-level semantic understanding in LVLMs.

# 5.2.3 ABLATIONS

In this section, we conduct a series of ablation studies under different experimental configurations to investigate which components contribute most to reducing hallucination in LVLMs, shown in Table 4. Specifically, we compare the following setups on the HalluText benchmark: (1) adding chain-of-thought (CoT) prompting, (2) simply appending raw OCR outputs to the prompt, and (3) our proposed OCRAssistor strategy. Results show that directly appending OCR results to the prompt brings modest gains on average. While CoT prompting yields some improvement, our OCRAssistor demonstrates substantially stronger gains, achieving improvements of 2.4% in localization, 13.1% in recognition, and 9.6% on the overall average metric. These results confirm that our carefully designed OCRAssistor effectively and seamlessly integrates OCR information into LVLMs. By leveraging the structured visual guidance provided by the OCR model, our method significantly alleviates OCR-aware hallucinations in both perception and understanding tasks.

# 5.2.4 OCR QUALITY

We also investigate the impact of OCR quality on hallucination mitigating, as shown in Table 5. Specifically, we evaluate the recognition quality of OCR models on the 1,500 original images used to construct the HalluText benchmark. In addition to our default OCR system PaddleOCR<sup>1</sup>, we compare with another widely used alternative, EasyOCR<sup>2</sup>. Experimental results indicate that EasyOCR achieves a 1-N.E.D. score that is 2.2 points lower than PaddleOCR, suggesting slightly inferior recognition performance. Correspondingly, under the same experimental settings, the downstream results on HalluText using EasyOCR are consistently lower than those with PaddleOCR. These findings demonstrate a positive correlation between OCR quality and hallucination mitigating performance: higher-quality OCR outputs provide more reliable visual cues, which better guide LVLMs and reduce hallucinated generations.

<sup>1</sup>https://github.com/PaddlePaddle/PaddleOCR

<sup>&</sup>lt;sup>2</sup>https://github.com/JaidedAI/EasyOCR

Table 4: Ablation for all components. Baseline selects Qwen2.5-VL-7B.

Model	$Acc_{loc}$	$Acc_{rec}$	$Acc_{all}$
Baseline	69.5	69.0	69.1
Baseline + CoT	70.5 (+1.0)	76.5 (+7.5)	74.6 (+5.6)
Baseline + OCR	70.1 (+0.6)	72.1 (+3.1)	71.4 (+2.4)
Baseline + OA	<b>71.9</b> (+2.4)	<b>82.1</b> (+13.1)	<b>78.7</b> (+9.6)

Table 5: Effect of different OCR models. 1-N.E.D. is a recognition metric defined as 1-NED.

OCR Model	1-N.E.D.	$Acc_{loc}$	$Acc_{rec}$	$Acc_{all}$
PaddleOCR	88.7	71.9	82.1	78.7
EasyOCR	86.5	71.5	79.1	76.6

Table 6: The efficiency experiments between OCRAssistor and VDGD (Ghosh et al., 2025). For convenience, we use Qwen2.5-VL-3B as the base model.

Settings	HalluText	Time(s/image)
Qwen2.5-VL-3B	59.3	0.275
Qwen2.5-VL-3B + VDGD	64.6	10.972
Qwen2.5-VL-3B + OA	71.4	0.817

Table 7: The gains of OCRAssistor on Qwen2.5-VL series across different scales.

Model	HalluText	OCRBench v2(EN/ZH)
Qwen2.5-VL-3B Qwen2.5-VL-7B	+12.1 +9.6	+1.3 / +4.7 +2.5 / +3.7
Qwell2.3-VL-7B	+9.0	+2.37 +3.7

# 5.2.5 SCALING

We further examine the effectiveness of OCRAssistor across different model scales, with results summarized in Table 7. Experimental results demonstrate that OCRAssistor consistently yields performance gains across both model sizes. Notably, the improvements are more pronounced on HalluText, which is explicitly designed to evaluate hallucination, indicating that our decoding strategy is particularly effective in hallucination-prone scenarios. These findings highlight the robust generalization ability of our method across LVLMs of varying capacity, making it applicable to both lightweight and large-scale models.

# 5.2.6 EFFICIENCY

Table 6 presents the runtime performance of OCRAssistor on Qwen2.5-VL. We compare three setups: the original 3B model, the VDGD-enhanced model, and our OCRAssistor. OCRAssistor introduces only 0.6s of additional latency per image while delivering a 12.1% performance gain. In contrast, VDGD adds over 10s per image, making it impractical despite modest improvements. This demonstrates OCRAssistor's favorable balance of efficiency and effectiveness. Notably, since our evaluation involves only multiple-choice outputs, baseline inference times remain low. In more complex scenarios requiring free-form or subjective generation, the overall inference latency would increase significantly, thereby reducing the relative overhead introduced by the OCR module. Thus, the efficiency advantage of our method could be more pronounced in real-world applications.

# 6 Conclusion

In this work, we present a comprehensive study on OCR-centric hallucinations in LVLMs. After applying a dual-perspective taxonomy that categorizes errors by task process (localization, recognition) and hallucination type (category, relation, attribute) and analyzing failure cases, we introduce HalluText, a fine-grained benchmark comprising 4,678 samples across 9 subsets, designed to diagnose OCR hallucinations. To address these challenges, we develop OCRAssistor, a training-free and plug-and-play pipeline that leverages external OCR signals to guide LVLM decoding. Experiments on HalluText and OCRBench v2 show that OCRAssistor consistently improves performance across models of different scales, while remaining efficient and scalable. Our findings underscore not only the importance of structured OCR integration but also highlight the effectiveness of a large-small model collaboration paradigm, where a lightweight OCR expert module supplements the strengths of a powerful LVLM. This cooperative design offers a practical and generalizable solution for reducing hallucinations in vision-language understanding.

# REPRODUCIBILITY STATEMENT

We have made extensive efforts to ensure that the results reported in this work are reproducible. All model architectures, training procedures, and hyperparameter settings are described in the main text (Sections 4–5) and detailed further in the Appendix (Appendix A–C). For the datasets used in our experiments, we provide complete descriptions of preprocessing and filtering steps in the supplementary materials. All evaluation metrics are formally defined in Section 3.3, enabling consistent replication of our analysis. Additionally, the source code and scripts used for training, inference, and evaluation will be made publicly available as anonymized supplementary material, facilitating direct reproduction of the reported results. Readers are referred to these resources for all necessary details to reproduce the experiments and analyses presented in this work.

# ETHICS STATEMENT

All authors have read and adhered to the ICLR Code of Ethics. This work focuses on analyzing and mitigating OCR hallucination, and does not involve direct experimentation on human subjects. All datasets used are either publicly available or used under appropriate licenses, and any personal information has been anonymized to protect privacy. We are aware of potential societal impacts of multimodal AI systems, including misuse for generating misleading content or biased outputs. In our experiments, we take care to evaluate model behavior across diverse languages and scenarios to mitigate unintended bias. No datasets or methods used are expected to cause harm to individuals or communities. We encourage responsible use and recommend that future users of the proposed models follow relevant legal, privacy, and fairness guidelines. Any conflicts of interest have been disclosed, and all research practices adhere to established standards of scientific integrity.

# REFERENCES

- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv:2502.13923*, 2025.
- Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou. Hallucination of multimodal large language models: A survey. *arXiv:2404.18930*, 2024.
- Ali Furkan Biten, Ruben Tito, Andres Mafla, Lluis Gomez, Marçal Rusinol, Ernest Valveny, CV Jawahar, and Dimosthenis Karatzas. Scene text visual question answering. In *ICCV*, pp. 4291–4301, 2019.
- Chee Kheng Ch'ng and Chee Seng Chan. Total-text: A comprehensive dataset for scene text detection and recognition. In *ICDAR*, volume 1, pp. 935–942. IEEE, 2017.
- Alessandro Favero, Luca Zancato, Matthew Trager, Siddharth Choudhary, Pramuditha Perera, Alessandro Achille, Ashwin Swaminathan, and Stefano Soatto. Multi-modal hallucination control by visual information grounding. In *CVPR*, pp. 14303–14312, 2024.
- Ling Fu, Zhebin Kuang, Jiajun Song, Mingxin Huang, Biao Yang, Yuzhe Li, Linghao Zhu, Qidi Luo, Xinyu Wang, Hao Lu, et al. Ocrbench v2: An improved benchmark for evaluating large multimodal models on visual text localization and reasoning. *arXiv:2501.00321*, 2024.
- Sreyan Ghosh, Chandra Kiran Reddy Evuru, Sonal Kumar, Utkarsh Tyagi, Oriol Nieto, Zeyu Jin, and Dinesh Manocha. Visual description grounding reduces hallucinations and boosts reasoning in lvlms. In *ICLR*, 2025.
- Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and preventing hallucinations in large vision language models. In *AAAI*, volume 38, pp. 18135–18143, 2024.
- Masato Hagiwara and Masato Mita. Github typo corpus: A large-scale multilingual dataset of misspellings and grammatical errors. *arXiv*:1911.12893, 2019.
- Haibin He, Maoyuan Ye, Jing Zhang, Xiantao Cai, Juhua Liu, Bo Du, and Dacheng Tao. Reasoning-ocr: Can large multimodal models solve complex logical reasoning problems from ocr cues? *arXiv*:2505.12766, 2025a.

- Zhentao He, Can Zhang, Ziheng Wu, Zhenghao Chen, Yufei Zhan, Yifan Li, Zhao Zhang, Xian Wang, and Minghui Qiu. Seeing is believing? mitigating ocr hallucinations in multimodal large language models, 2025b. URL https://arxiv.org/abs/2506.20168.
  - Mingxin Huang, Yuliang Liu, Dingkang Liang, Lianwen Jin, and Xiang Bai. Mini-monkey: Alleviate the sawtooth effect by multi-scale adaptive cropping. *arXiv e-prints*, pp. arXiv–2408, 2024a.
  - Mingxin Huang, Yongxin Shi, Dezhi Peng, Songxuan Lai, Zecheng Xie, and Lianwen Jin. Ocrreasoning benchmark: Unveiling the true capabilities of mllms in complex text-rich image reasoning. *arXiv*:2505.17163, 2025.
  - 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. In *CVPR*, pp. 13418–13427, 2024b.
  - Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and CV Jawahar. Icdar2019 competition on scanned receipt ocr and information extraction. In *ICDAR*, pp. 1516–1520. IEEE, 2019.
  - Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. arXiv:2410.21276, 2024.
  - Jitesh Jain, Jianwei Yang, and Humphrey Shi. Vcoder: Versatile vision encoders for multimodal large language models. In *CVPR*, pp. 27992–28002, 2024.
  - Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Funsd: A dataset for form understanding in noisy scanned documents. In *ICDARW*, volume 2, pp. 1–6. IEEE, 2019.
  - Qing Jiang, Jiapeng Wang, Dezhi Peng, Chongyu Liu, and Lianwen Jin. Revisiting scene text recognition: A data perspective. In *ICCV*, pp. 20543–20554, 2023.
  - Dimosthenis Karatzas, Faisal Shafait, Seiichi Uchida, Masakazu Iwamura, Lluis Gomez i Bigorda, Sergi Robles Mestre, Joan Mas, David Fernandez Mota, Jon Almazan Almazan, and Lluis Pere De Las Heras. Icdar 2013 robust reading competition. In *ICDAR*, pp. 1484–1493. IEEE, 2013.
  - Solomon Kullback. Kullback-leibler divergence. Tech. Rep., 1951.
  - 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. In *CVPR*, pp. 13872–13882, 2024.
  - Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv:2407.07895*, 2024a.
  - Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. Monkey: Image resolution and text label are important things for large multi-modal models. In *CVPR*, pp. 26763–26773, 2024b.
  - Zhang Li, Yuliang Liu, Qiang Liu, Zhiyin Ma, Ziyang Zhang, Shuo Zhang, Zidun Guo, Jiarui Zhang, Xinyu Wang, and Xiang Bai. Monkeyocr: Document parsing with a structure-recognition-relation triplet paradigm. *arXiv*:2506.05218, 2025.
  - Chongyu Liu, Yuliang Liu, lianwen Jin, Shuaitao Zhang, Canjie Luo, and Yongpan Wang. Erasenet: End-to-end text removal in the wild. *IEEE TIP*, 29:8760–8775, 2020.
    - Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Mitigating hallucination in large multi-modal models via robust instruction tuning. *arXiv*:2306.14565, 2023.
    - Yuliang Liu, Lianwen Jin, Shuaitao Zhang, and Sheng Zhang. Detecting curve text in the wild: New dataset and new solution. *arXiv*, 2017.

- Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. *SCIS*, 67(12):220102, 2024.
  - Shiyin Lu, Yang Li, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, and Han-Jia Ye. Ovis: Structural embedding alignment for multimodal large language model. *arXiv*:2405.20797, 2024.
  - Colin M MacLeod. Half a century of research on the stroop effect: an integrative review. *Psychological bulletin*, 109(2):163, 1991.
  - Fiona Macpherson and Dimitris Platchias. *Hallucination: Philosophy and psychology*. MIT Press, 2013.
  - Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv:2203.10244*, 2022.
  - Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. Infographicvqa. In *WACV*, pp. 1697–1706, 2022.
  - Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and factuality in abstractive summarization. *arXiv:2005.00661*, 2020.
  - Mor Shpigel Nacson, Aviad Aberdam, Roy Ganz, Elad Ben Avraham, Alona Golts, Yair Kittenplon, Shai Mazor, and Ron Litman. Docvlm: Make your vlm an efficient reader. In *CVPR*, pp. 29005–29015, 2025.
  - Linke Ouyang, Yuan Qu, Hongbin Zhou, Jiawei Zhu, Rui Zhang, Qunshu Lin, Bin Wang, Zhiyuan Zhao, Man Jiang, Xiaomeng Zhao, et al. Omnidocbench: Benchmarking diverse pdf document parsing with comprehensive annotations. In *CVPR*, pp. 24838–24848, 2025.
  - Yan Shu, Hangui Lin, Yexin Liu, Yan Zhang, Gangyan Zeng, Yan Li, Yu Zhou, Ser-Nam Lim, Harry Yang, and Nicu Sebe. When semantics mislead vision: Mitigating large multimodal models hallucinations in scene text spotting and understanding. *arXiv*:2506.05551, 2025.
  - Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *CVPR*, pp. 8317–8326, 2019.
  - Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. *arXiv:2309.14525*, 2023.
  - Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv:2403.05530*, 2024.
  - Chenxi Wang, Xiang Chen, Ningyu Zhang, Bozhong Tian, Haoming Xu, Shumin Deng, and Huajun Chen. Mllm can see? dynamic correction decoding for hallucination mitigation. *arXiv:2410.11779*, 2024a.
  - Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv:2409.12191*, 2024b.
  - Xintong Wang, Jingheng Pan, Liang Ding, and Chris Biemann. Mitigating hallucinations in large vision-language models with instruction contrastive decoding. *arXiv:2403.18715*, 2024c.
  - Haoran Wei, Chenglong Liu, Jinyue Chen, Jia Wang, Lingyu Kong, Yanming Xu, Zheng Ge, Liang Zhao, Jianjian Sun, Yuang Peng, et al. General ocr theory: Towards ocr-2.0 via a unified end-to-end model. *arXiv preprint arXiv:2409.01704*, 2024.
  - Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. arXiv:2408.01800, 2024.

- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *NSR*, 11(12):nwae403, 2024.
  - Qifan Yu, Juncheng Li, Longhui Wei, Liang Pang, Wentao Ye, Bosheng Qin, Siliang Tang, Qi Tian, and Yueting Zhuang. Hallucidoctor: Mitigating hallucinatory toxicity in visual instruction data. In *CVPR*, pp. 12944–12953, 2024a.
  - Ya-Qi Yu, Minghui Liao, Jihao Wu, Yongxin Liao, Xiaoyu Zheng, and Wei Zeng. Texthawk: Exploring efficient fine-grained perception of multimodal large language models. *arXiv:2404.09204*, 2024b.
  - Jinrui Zhang, Teng Wang, Haigang Zhang, Ping Lu, and Feng Zheng. Reflective instruction tuning: Mitigating hallucinations in large vision-language models. In *ECCV*, pp. 196–213. Springer, 2024.
  - Zhen Zhao, Jingqun Tang, Binghong Wu, Chunhui Lin, Shu Wei, Hao Liu, Xin Tan, Zhizhong Zhang, Can Huang, and Yuan Xie. Harmonizing visual text comprehension and generation. *arXiv:2407.16364*, 2024.
  - Zhiyuan Zhao, Bin Wang, Linke Ouyang, Xiaoyi Dong, Jiaqi Wang, and Conghui He. Beyond hallucinations: Enhancing lvlms through hallucination-aware direct preference optimization. *arXiv*:2311.16839, 2023.
  - Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv:2504.10479*, 2025a.
  - Lanyun Zhu, Deyi Ji, Tianrun Chen, Peng Xu, Jieping Ye, and Jun Liu. Ibd: Alleviating hallucinations in large vision-language models via image-biased decoding. In *CVPR*, pp. 1624–1633, 2025b.

The appendix includes the following aspects:

- A: Use of Large Language Models
- B: Details of HalluText curation.
- C: Details of prompt.
- D: Additional Experiments.
- E: Visualization.

702

703 704

705

706

708

709

710 711

712 713

714

715

716

717

718

719

720 721

722723724

750 751 752

753

754

755

# A USE OF LARGE LANGUAGE MODELS

In this work, large language models (LLMs) are used solely as generally purpose assistive tools to improve the clarity, grammar, and readability of the manuscript. LLMs are not used for research ideation, data analysis, model development, or any other scientific decision-making. All scientific content, ideas, results, and conclusions presented in this paper are independently produced by the authors. The authors take full responsibility for the accuracy and integrity of the work, including any content that was refined or edited with the assistance of LLMs. No information generated by LLMs that could constitute plagiarism, fabrication, or scientific misconduct has been included.

# B DETAILS OF HALLUTEXT CURATION

# Algorithm 1 Generating Stroop-Effect QA Pairs

```
725
       Require: colors, shapes
726
       Ensure: qa_pair
727
        1: function GENERATESTROOPQA
728
        2:
              shape_color, text_color, render_color ← RandomSelect(colors)
        3:
              shape_shape, text_shape ← RandomSelect(shapes)
729
        4:
              img \( \text{ImageDraw}(\text{shape_shape, shape_color, render_color) \)
730
        5:
              text ← concatenate(text_color, text_shape)
731
              font\_size \leftarrow RandomSelect(range(10, 50))
        6:
732
        7:
              bbox \leftarrow ComputeBBox(text, font\_size)
733
        8:
              if apply_rotation then
734
        9:
                  rotated_dims ← GetRotatedDims(bbox, angle)
735
       10:
                  if rotated_dims exceeds image boundaries then
736
       11:
                     font_size ← AdjustFontSize(font_size, max_scale)
737
                     bbox \leftarrow ComputeBBox(text, font\_size)
       12:
738
       13:
                  end if
739
                  pos ← FindValidPos(rotated_dims)
       14:
       15:
                  RenderTextRotated(text, pos, font_size, angle)
740
       16:
              else
741
                  pos \leftarrow FindValidPos(bbox)
       17:
742
       18:
                  RenderText(text, pos, font_size)
743
       19:
744
              question ← "What text is written in the image?"
       20:
745
              options ← GenerateOptions(text_color,
                                                                                shape_color,
                                                                text_shape,
746
           shape_shape, render_color)
747
       22.
              options, answer ← ShuffleOptions(text, optiobs)
748
       23:
              qa_pair ← {question, img, options, answer}
749
       24: end function
```

## B.1 EXISTENCE

We construct the Existence subset using SCUT-ENS Liu et al. (2020), a dataset containing paired images before and after scene text erasing. By leveraging these image pairs and their corresponding OCR annotations, we create VQA-style samples that ask whether a specific text instance existed

Table 8: Prompt templates for different settings.

Settings	Prompt
	HalluText
Baseline	Please strictly follow these rules: \n Place the answer only option letter (with no extra characters) on a separate last line. \n Question: [QUESTION]. \n Options: [OPTIONS]. \n Answer: \n
Baseline+CoT	Please strictly follow these rules: \n Let us think this question step by step (Chain of thought) and Place the answer only option letter (with no extra characters) on a separate last line. \n Question: [QUESTION]. \n Options: [OPTIONS]. \n Chain of thought: \n Answer: \n
Baseline+OCR	The texts in image was recognized in the image: [OCR RESULTS] \Please strictly follow these rules: \n Place the answer only option letter (with no extra characters) on a separate last line. \n Question: [QUESTION]. \n Options: [OPTIONS].\n Answer: \n
Baseline+OCRAssistor	The texts in image was recognized in the image: [OCR RESULTS]. Please strictly follow these rules: \n Let us think this question step by step (Chain of thought) and Place the answer only option letter (with no extra characters) on a separate last line. \n Question: [QUESTION]. \n Options: [OPTIONS]. \n Chain of thought: \n Answer: \n

prior to erasure. As shown in Figure 5, we design the questions template "Does the text 'TEA' exist in the image?", accompanied by the erased image as shown in Figure 5 (a). The correct answer is clearly No. If the LVLM relies solely on dataset bias rather than visual information provided by the user, it is prone to incorrectly predicting Yes. To further balance the distribution of answers, we deliberately incorporate negative forms in the question design, ensuring a more even ratio between Yes and No responses.

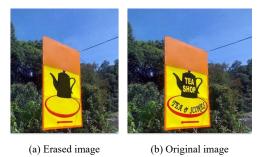


Figure 5: The details of "Existence" subset.

# B.2 INCOMPLETE

The *Incompletion* subset is adapted from OCRBench v2 using its standard recognition prompts. To ensure quality, we manually verify and clean the original annotations. For each question, a confusion option is generated by using the original word, which is rich in semantics. The other two distractors are created by applying random character-level edits (insertion, deletion, substitution) based on the ground truth. An example is shown in Figure 3.

# B.3 COUNTING

The *Counting* subset also follows the instruction of OCRBench v2, using standard counting-style prompts. Answers are derived from the original dataset, with confusion options introduced by sampling positive integer near the correct answer (e.g. +1, -1) to simulate realistic ambiguity.

# B.4 TYPO

The *Typo* subset is synthetically constructed using the typo corpus and the Pillow image library. We adopt a minimalist rendering black text on a white background, without any decorative elements. Confusion options are generated by randomly applying one character-level edit (insertion, deletion, or substitution) to the ground truth or its corrected version. This process mirrors that of the Incomplete subset, focusing on recognition robustness under typographical noise.

# B.5 Position

The *Position* subset is built from the Total-Text dataset. We discard illegible or unannotated text instances and classify the remaining ones into eight relative positional categories: top-left, top, top-right, right, bottom-right, bottom, bottom-left, and left. To avoid ambiguity between adjacent classes (e.g., top-left vs. top), we explicitly remove potentially confusing categories when generating answer options. This ensures that each question has one unambiguous correct answer.

# B.6 INDEX, SLICE, AND FREQUENCY

The *Index*, *Slice*, and *Frequency* subsets are jointly derived from the Union14M-contextless dataset Jiang et al. (2023). For each image–annotation pair, we sample character-level statistics such as frequency, position, and index to formulate distinct question types. To maintain a single-answer format, we filter out ambiguous cases—such as words with repeated characters—where multiple valid answers might exist for an index-based query.

#### B.7 STROOP EFFECT

The *Stroop Effect* subset is uniquely constructed without relying on any existing public dataset. We manually generate image—question—answer triplets to simulate conditions where irrelevant but plausible distractors interfere with OCR perception. The generation pipeline is detailed in Algorithm 1. Crucially, all confusion options used in the answer choices are explicitly present within the image, enabling a faithful evaluation of the LVLM's susceptibility to OCR hallucinations.

# C DETAILS OF PROMPT

To adapt different ablation settings, we design prompt templates tailored to various input configurations. The detailed prompt formats are provided in Table 8. In the HalluText benchmark, the *Baseline* configuration includes only the core question and the corresponding answer choices. For the *CoT* and *OCR* settings, we incorporate respective guiding cues into the prompt. In the OCRAssistor setup, we include both the CoT prompt and OCR information in the language instruction to encourage alignment between the LVLM's output and the OCR-derived content distribution. Experimental results show that our method significantly improves performance on HalluText and has a consistent positive effect in mitigating OCR hallucinations. On the OCRBench v2 benchmark, we follow the standard evaluation protocol, using only the question as the full prompt. Under the fair setting, our method demonstrates stable and consistent improvements in fair comparisons with other models, as shown in Table 3.

Table 9: Ablation for OCR inputs. The baseline LVLM is Qwen2.5-VL-3B. All settings loads OCRAssistor.

Model	$Acc_{loc}$	$Acc_{rec}$	$Acc_{all}$
Rec-only	61.7	76.3	71.4
Det & Rec	63.5 (+1.8)	72.8 (-2.5)	69.7 (-1.7)

Table 10: Ablation for the setting of  $\lambda$ , The baseline LVLM is Qwen2.5-VL-3B. **Bold** indicates the best performance.

λ	$Acc_{loc}$	$Acc_{rec}$	$Acc_{all}$
0.1	60.3	76.1	70.8
0.5	61.7	76.3	71.4
1.0	60.1	76.2	70.8
1.5	60.5	75.9	70.7
2.0	60.5	76.2	70.9

Table 11: Detailed results of Qwen2.5-VL-3B on OCRBench v2. OA indicates the OCRAssistor. Abbreviations: TR = Text Recognition, TD = Text Detection, TS = Text Spotting, RE = Relation Extraction, EP=Element Parsing, MC = Metathetical Calculating, TU=Text Understanding, KR = Knowledge Reasoning.

Model	English Part								Chinese Part						Overall	
Model	TR	TD	TS	RE	EP	MC	TU	KR	TR	RE	EP	TU	KR	English	Chinese	
Qwen2.5-VL-3B	63.9	18.7	0.0	81.5	32.5	35.3	69.2	49.2	69.0	47.2	33.0	35.5	43.5	43.8	45.6	
Qwen2.5-VL-3B+OA	58.9	20.6	0.0	84.7	34.6	39.9	70.9	51.0	67.6	54.8	33.2	54.0	41.8	45.1 (+1.3)	50.3 (+4.7)	

# D ADDITIONAL EXPERIMENTS

In this section, we provide additional experimental results that are omitted from the main text due to space limitations.

# D.1 OCR INPUTS

We conduct an ablation study on the use of OCR inputs. Two configurations are compared: (1) using only the OCR recognition results as input, and (2) incorporating both detection and recognition results into the prompt. As shown in Table 9, providing both detection and recognition results as OCR priors leads to a 1.8% improvement on the localization task compared to using recognition results alone. However, this setting results in performance drops of 2.5% and 1.7% on the recognition task and the overall average, respectively. We attribute this phenomenon to the limited guidance provided by the coordinate-format detection results after tokenization, which could not be effectively utilized during LVLM decoding. We caution that directly including OCR detection outputs in the prompt make adverse effects.

# D.2 The effect of $\lambda$

 $\lambda$  is the guidance factor used in OCRAssistor sampling, as defined in Equation (5). A smaller  $\lambda$  indicates weaker influence. We evaluate the impact of  $\lambda$  on hallucination mitigation using the HalluText benchmark. Table 10 reports results on the Qwen2.5-VL-3B model. The experiments show that the best performance is achieved when  $\lambda=0.5$ . Moreover, the model exhibits relatively stable average performance when  $\lambda$  is within the range of 0.1 to 2.

# 

# D.3 THE DETAILED RESULTS ON QWEN2.5-VL-3B

Owing to space constraints, Table 7 presents only the performance gains of Qwen2.5-VL-3B with OCRAssistor on HalluText and OCRBench v2. For completeness, the detailed results are provided in Table 11 and Table 12.

# E VISUALIZATION

This section presents qualitative visualizations of Qwen2.5-VL-7B's performance on two datasets.

Table 12: Detailed performance of Qwen2.5VL-3B on HalluText. OA indicates the OCRAssistor. Abbreviations: EX = Existence, RP = Relative Position, CT = Counting, ST = Stroop Effect, TY = Typo, IC = Incompletion, FQ = Frequency, ID = Index.

Model	Localization			4			1	1.00				
Model	EX	POS	CT	$Acc_{loc}$	ST	TY	IC	FQ	ID	SL	$Acc_{rec}$	$Acc_{all}$
Qwen2.5-VL-3B	90.8	39.4	46.4	58.9	94.0	79.2	42.1	49.5	45.0	47.7	59.6	59.3
Owen2.5-VL-3B + OA	91.6	41.0	52.4	61.7	100.0	80.0	64.3	67.8	68.6	77.1	76.3	71.4

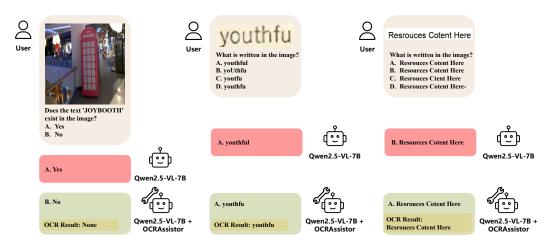


Figure 6: A Visualization of HalluText.

# E.1 HALLUTEXT

Figure 6 illustrates results on HalluText, where we visualize the input image, question, answer options, model predictions (before and after enhancement), and the OCR-recognized text. The comparison shows that with guidance from OCR outputs, Qwen2.5-VL-7B better adheres to visual instructions and exhibits reduced OCR hallucinations.

# E.2 OCRBENCH V2

Figure 7 and Figure 8 show visualizations on the English and Chinese subsets of OCRBench, respectively. We observe that the proposed OCRAssistor module helps the LVLM correct fine-grained recognition errors. For example, in Figure 7, the model originally extracted "Newspaper Parent" for the field "Brand(s) Applicable", while the image text actually reads "Newport Parent"; similarly, it misread "Coupon Issue Date" as "4/1/00" instead of the correct "4/14/00". These cases highlight the presence of OCR hallucinations in the baseline LVLM, which are significantly mitigated after applying the proposed improvements. In summary, our method achieves stable performance gains across diverse generalized OCR scenarios.

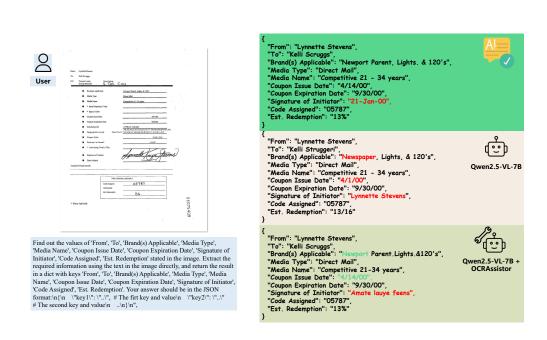


Figure 7: A Visualization of OCRBench v2-EN.



Figure 8: A Visualization of OCRBench v2-CN.