# OCRBench v2: An Improved Benchmark for Evaluating Large Multimodal Models on Visual Text Localization and Reasoning

**Ling Fu**[1]   **Zhebin Kuang**[1]   **Jiajun Song**[1]   **Mingxin Huang**[2]   **Biao Yang**[1]
**Yuzhe Li**[1]   **Linghao Zhu**[1]   **Qidi Luo**[1]   **Xinyu Wang**[3]   **Hao Lu**[1]   **Zhang Li**[1]
**Guozhi Tang**[4]   **Bin Shan**[4]   **Chunhui Lin**[4]   **Qi Liu**[4]   **Binghong Wu**[4]
**Hao Feng**[4]   **Hao Liu**[4]   **Can Huang**[4]   **Jingqun Tang**[4]   **Wei Chen**[1]
**Lianwen Jin**[2]   **Yuliang Liu**[1][*]   **Xiang Bai**[1][*]

[1]Huazhong University of Science and Technology   [2]South China University of Technology
[3]University of Adelaide   [4]ByteDance

## Abstract

Scoring the Optical Character Recognition (OCR) capabilities of Large Multimodal Models (LMMs) has witnessed growing interest. Existing benchmarks have highlighted the impressive performance of LMMs in text recognition; however, their abilities in certain challenging tasks, such as text localization, handwritten content extraction, and logical reasoning, remain underexplored. To bridge this gap, we introduce **OCRBench v2**, a large-scale bilingual text-centric benchmark with currently the most comprehensive set of tasks ($4\times$ more tasks than the previous multi-scene benchmark OCRBench), the widest coverage of scenarios (31 diverse scenarios), and thorough evaluation metrics, with $10,000$ human-verified question-answering pairs and a high proportion of difficult samples. Moreover, we construct a private test set with $1,500$ manually annotated images. The consistent evaluation trends observed across both public and private test sets validate the OCRBench v2's reliability. After carefully benchmarking state-of-the-art LMMs, we find that most LMMs score below 50 (100 in total) and suffer from five-type limitations, including less frequently encountered text recognition, fine-grained perception, layout perception, complex element parsing, and logical reasoning. The benchmark and evaluation scripts are available at https://github.com/Yuliang-Liu/MultimodalOCR.

## 1   Introduction

The emergence of Large Language Models (LLMs) [1, 2, 3] has greatly improved the understanding and generation of structured text. However, in reality, much of the textual content is unstructured; it appears within images, videos, and other non-textual media in varied positions, orientations, and shapes. The need for processing such unstructured content leads to the study of Large Multimodal Models (LMMs) [4, 5, 6] that extend the text-only LLMs to additional modalities. By pretraining on multimodal data, LMMs acquire the zero-shot ability to interpret across diverse media, such as recognizing and understanding complex visual scene text [7]. Such capability represents a significant advancement over standard Optical Character Recognition (OCR), because LMMs not only spot text but also interpret its semantic relevance to a scene.

Compared with classic OCR that typically relies on task-specific models to spot text, the increasing capability of LMMs to process multimodal inputs has opened new potential to redefine the area of
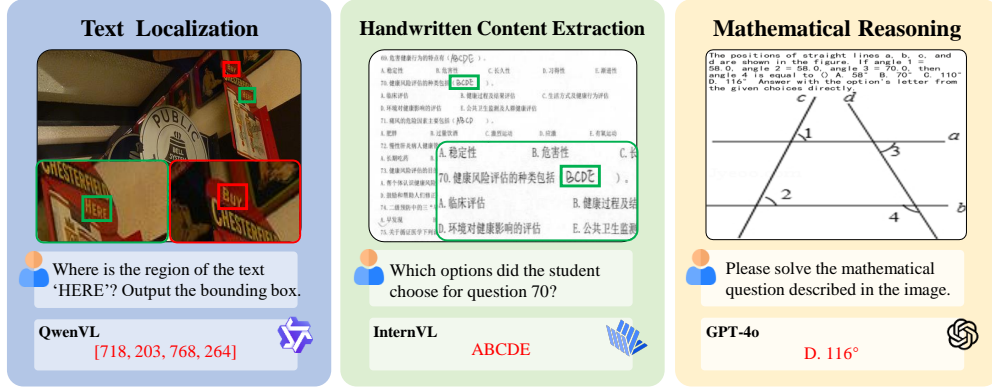
---

[*]Corresponding authors

Figure 1: **Large multimodal models struggle with text-intensive tasks accurately**. They are prone to errors in tasks like text localization, handwritten content extraction, and mathematical reasoning, revealing limitations in tackling complex textual information within images.

OCR. OCR has therefore become an important aspect of recent LMM evaluations. Some text-focused tasks have been included in standard benchmarks to assess the proficiency of LMMs in recognizing and interpreting textual content [8, 9]. Typically, text-based Visual Question Answering (VQA) datasets [10, 11, 12] are repurposed to evaluate OCR by framing generic VQA into questions that require accurate reading of embedded text. However, many of these datasets are initially created for classic OCR models, which are of limited diversity, depth, and suitability for evaluating LMMs. A common drawback is that, many questions lack sufficient complexity to assess the reasoning abilities of LMMs on scene text, and some can even be answered without visual input [13, 12].

More recently, several customized benchmarks [14, 15, 16, 17, 18] have explored the OCR capabilities of LMMs. For example, OCRBench [14] consolidates 5 core text-oriented tasks to evaluate LMM performance across traditional OCR functions. Other datasets, such as ComTQA [19] and ChartX [20], focus on structured text interpretation like table and chart understanding. While such effort represents a leap over standard OCR benchmarks, they remain limited in both data diversity and quantity (see Tab. 1), often leading to rapid performance saturation. For example, recent LMMs such as Qwen2.5-VL [21] have achieved 96.4% accuracy on the DocVQA dataset [22], nearly matching human performance at 98.1%, and 88.8% on OCRBench [14]. This raises an important question for the community: *Do models perform well enough on text-oriented visual understanding tasks in the LMM era, or do existing benchmarks fail to capture the broader challenges in diverse environments?*

To answer the question above, we conducted preliminary tests with several state-of-the-art LMMs, including Qwen2.5-VL-7B [21], InternVL3-14B [23], and GPT-4o [24]. These tests assessed performance on text-oriented tasks, such as text localization, handwritten content extraction, and document-based logical reasoning. As illustrated in Fig. 1, each model can fail on one of the text-intensive tasks. These failures reveal a gap in detailed visual perception across different models, which constrains their effectiveness in tasks requiring accurate text localization, recognition, and contextual understanding within images. Recent benchmarks, such as OmniDocBench [25], CC-OCR [26], and MMLONGBENCH-DOC [27], have broadened evaluation to cover more comprehensive scenarios, including fine-grained document parsing and multi-page document understanding. Their analyses reveal the limited capabilities of LMMs for practical OCR applications and highlight the growing need for benchmarks that allow for more robust and varied evaluation of LMMs.

To bridge this gap, we propose *OCRBench v2*, a comprehensive benchmark designed to assess LMMs across diverse text-oriented visual understanding tasks. As shown in Fig. 3, *OCRBench v2* assesses eight core text-reading abilities, including *text recognition*, *text referring*, *text spotting*, *relation extraction*, *element parsing*, *mathematical calculation*, *visual text understanding*, and *knowledge reasoning*, organized into a total of 23 concrete tasks. This benchmark provides $10,000$ high-quality, human-validated instruction-response pairs and also six types of evaluation metrics, which offers a rigorous framework for evaluating LMM performance in complex, practical OCR scenarios. For better evaluation quality, we further collect and label $1,500$ additional text-images from scratch, reserved as the private test set. This private data serves as an independently curated test set to validate model generalization. In summary, the contributions of this work are three-fold:

Table 1: Comparison between the proposed benchmark and existing text-centric datasets.

| Benchmark | #Scenario | #Task | #Image | #Instruction |
|---|---|---|---|---|
| OCRbench [14] | $\sim 14$ | 5 | 0.9k | 1k |
| Seed-bench-2-plus [15] | $\sim 8$ | 1 | 0.6k | 2.3k |
| CONTEXTUAL [16] | $\sim 11$ | 1 | 0.5k | 0.5k |
| Fox [17] | 2 | 9 | 0.7k | 2.2k |
| MMTab-eval [28] | 1 | 9 | 23k | 49k |
| ComTQA [19] | 1 | 4 | 1.6k | 9k |
| ChartX [20] | 1 | 7 | 6k | 6k |
| MMC [29] | 1 | 9 | 1.7k | 2.9k |
| OmniDocBench [25] | 9 | 5 | 1k | 1k |
| MMLONGBENCH-DOC [27] | 7 | 2 | 6.4k | 1.1k |
| OCRBench v2 (Ours) | 31 | 23 | 9.5k | 10k |

- *OCRBench v2*: an improved benchmark designed to assess eight core OCR competencies and covers 23 tasks across 31 diverse scenarios, which provides a thorough evaluation framework encapsulating fundamental and advanced text-centric challenges.

- We systematically evaluate state-of-the-art LMMs, ranging from commercial APIs to open-source models, which establishes broad baselines for OCR performance and enables a comparative understanding of model capabilities across varied text-oriented visual understanding tasks.

- We provide a detailed analysis to identify factors affecting the OCR capabilities of LMMs. The analysis examines performance across various dimensions such as model generalization to diverse text types, model robustness, and the ability to tackle complex visual-textual relations.

## 2 Related Work

**OCR-Enhanced LMMs.** Inspired by LLMs, visual encoders are integrated into them to create LMMs capable of processing both images and text. Early LMMs exhibit strong zero-shot OCR capabilities, motivating the exploration of text-centric LMMs. For instance, some work [30, 31] use text-centric instruction-tuning to enhance OCR-related abilities. But they are restricted to low-res inputs, limiting the ability to recognize dense and small text. To address this, several studies [32, 33, 34] shift attention to increasing the input resolution. As the resolution of inputs increases, so does computational cost. To tackle this issue, TextMonkey [7] introduces a Token Resampler to compress redundant visual feature tokens, mPLUG-DocOwl2 [35] presents a DocCompressor module for compressing high-res images, and DocKylin [36] adopts adaptive pixel slimming and dynamic token slimming modules to reduce redundant regions. To enhance layout perception, DocLayLLM [37] integrates layout information into LMMs inputs, LayTokenLLM [38] shares position IDs between text and layout tokens, DocMark [39] utilizes adaptive generation of markup languages to build structured document representations, while Marten [40] introduces an additional mask generator during pre-training. Despite strong results on existing benchmarks, challenges remain unsolved in certain key areas such as text localization, entity extraction, and logical reasoning.

**Benchmarks for Text-Centric LMMs.** Previous efforts have focused on creating scenario-specific benchmarks to assess LMMs. For example, DocVQA [22], ChartQA [41], Infographics VQA [42], and TextVQA [10] evaluate models on document understanding, chart reasoning, infographic interpretation, and scene text comprehension, respectively. To broaden evaluation scope, OCRBench [14] introduces a holistic evaluation framework covering five text-oriented tasks, while CONTEXTUAL [16] and SEED-Bench-2-Plus [15] introduce context-sensitive and diverse real-world images. Other benchmarks target specific challenges such as dense text understanding [43], complex structure parsing [26], and fine-grained document analysis [25]. To provide a more thorough assessment, some benchmarks design multiple tasks within a specific scenario. TableVQA-Bench [18], MMTab [28], and ComTQA [19] explore table-based tasks, while ChartY [44], ChartX [20], and MMC [29] focus on chart information extraction and reasoning. OmniDocBench [25] focuses on document parsing tasks and provides a comprehensive evaluation framework. Recently, DUDE [45], MM-NIAH [46], MP-DocVQA [47], MMLONGBENCH-DOC [27], and LongDocURL [48] explore the long document understanding capability of LMMs. In this work, we establish *OCRBench v2*, a systematic benchmark to reveal the limitations of LMMs in diverse single-image, text-related scenarios.
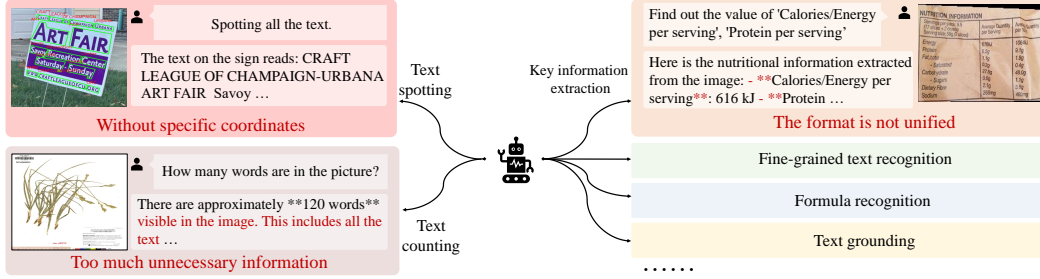
Figure 2: As evaluation for LMMs expands to diverse text-oriented tasks, existing datasets often require task-specific handling, making unified and scalable evaluation difficult.

## 3 Why Do We Need OCRBench v2?

**Limitations of Existing Benchmarks.** Recent evaluations of LMMs' OCR capabilities have made significant progress, yet most existing benchmarks exhibit limitations. Datasets like DocVQA, ChartQA, and TextVQA are often narrow in scope, focusing predominantly on text recognition within specific domains such as forms, tables, or documents. While useful for isolated capabilities, they fall short in task diversity, instruction complexity, and structured output formats that better reflect the multimodal nature of LMMs. In particular, many of these benchmarks were originally tailored for traditional OCR systems that prior to the emergence of LMMs. Furthermore, as illustrated in Fig. 2, complex task-specific processes are needed for LMMs when extended to more text-oriented tasks, which limits the evaluation of their broader capabilities. In this spirit, OCRBench v2 aims to evaluate OCR systems in terms of what ultimately matters: can a model recognize, understand, and reason over visual text to produce correct and meaningful answers?

**The Necessity of Unified Multi-task Evaluation.** With the emergence of LMMs, current models now excel at end-to-end performance across diverse tasks. Therefore, modern OCR goes beyond basic character recognition. Real-world documents often involve complex layouts and semantic structures that demand contextual understanding and reasoning. To assess these multi-task models, unified benchmarks like LongDocURL [48], OmniDocBench [25], CCOCR [26], OCRBench [14], CON-TEXTUAL [16], SEED-Bench-2-Plus [15], have been proposed and successfully demonstrated the value of evaluating text-oriented models across diverse tasks. These benchmarks show the importance of unified evaluation frameworks in guiding model development. However, as model capabilities expand, existing benchmarks with limited task coverage result in fragmented and sometimes misleading insights. To address this, a unified benchmark is essential to: 1) *Understand generalization*: Can a model perform consistently across varied text-centric tasks? 2) *Diagnose failure models*: Does a model that excels in recognition also succeed in reasoning, localization, and parsing? 3) *Guide model development*: Unified evaluation provides clearer signals for architecture and training improvements.

As shown in Fig. 3, *OCRBench v2* tackles this by combining 23 tasks under 8 core capabilities within one framework. This holistic design enables systematic comparison of models and highlights trade-offs (e.g., performance on reasoning vs. recognition) that isolated benchmarks cannot reveal.

**How OCRBench v2 Addresses the Gaps.** *OCRBench v2* is a comprehensive, and high-difficulty benchmark specifically built to evaluate LMMs in realistic OCR settings, with key advantages: 1) *Breadth of coverage*: With 31 scenarios, we ensure diverse contextual challenges; 2) *Task variety*: The benchmark spans 8 OCR-related capabilities, many of which are poorly handled by current LMMs; 3) *Instruction complexity*: Human-authored prompts and structured outputs (e.g., Markdown, JSON, LaTeX) raise the bar beyond simple answer extraction; 4) *Private evaluation test set*: To prevent overfitting and training contamination, we additionally provide a private test set.

Ultimately, *OCRBench v2* fills a critical gap by offering a unified and challenging benchmark that reflects the practical needs of OCR in the LMM era. It not only measures what current models can do, but more importantly, reveals what they still cannot.

**Design Rationale: Focusing on Single-Image Text Tasks.** While designing *OCRBench v2*, we focus on challenges in single-image, text-related scenarios, and do not extend our study to multi-image tasks. This design choice is grounded in two considerations: 1) Single-image understanding is the
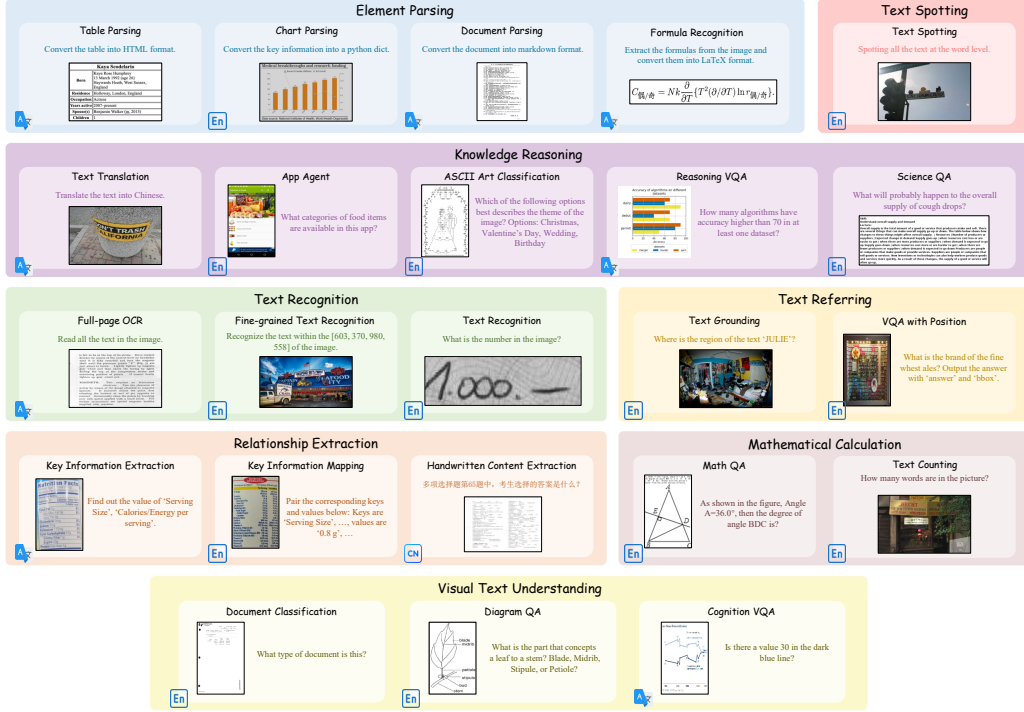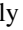
Figure 3: **Sample visualizations for each task.** OCRBench v2 comprises 23 sub-tasks grouped under 8 core OCR capabilities. Tasks marked with [A文] contain both English and Chinese instructions, while other tasks are either English-only [En] or Chinese-only [CN] (Zoomed in for better clarity).

foundation for more complex multimodal tasks. Many existing models still perform unsatisfactorily in various single-image scenarios, which motivates our work; 2) Given long-context inputs, multi-page tasks have more emphasis on long-sequence modeling, requiring specific benchmarks to assess this capability individually. For example, MMLONGBENCH-DOC focuses on evaluating the ability of LMMs to locate and understand content across pages in long documents.

**Private Dataset for Reliable Evaluation.** To further enhance the assessment quality, we also construct a private test set. This data comprises $1,500$ manually collected text-rich images with human-annotated labels, covering 23 tasks aligned with the distribution of the public data. Among the private data, 735 images were manually captured, and 765 images were sourced from unlabeled data with diverse scenarios. The data sources include printed books, e-books, scanned documents, and web content. During data collection and annotation, we meticulously curated samples to align with practical text-oriented applications. Given that benchmarks may be contaminated in massive internet-scraped pre-training data of LMMs, this data will not be released. Instead, we maintain a regularly updated leaderboard to reflect the performance of advanced LMMs. Moreover, consistent performance trends and model rankings observed on both the public and private test sets (see Section 5.2) indicate the benchmark's well-founded design and its effectiveness in identifying model capabilities.

# 4 Benchmark Construction

In this section, we describe the task description, annotation curation, statistics, and evaluation criteria.

## 4.1 Task Description

To provide a comprehensive evaluation framework for text-reading tasks, we categorize OCR capabilities into eight core areas, each encompassing specific sub-tasks that address various aspects of

text comprehension and interpretation. Fig. 3 exhibits samples for each task, with visual inputs and corresponding instructions. Detailed descriptions of these core capabilities are as follows.

**Text Recognition.** This fundamental capability focuses on perceiving textual content. The related tasks include (fine-grained) text recognition and full-page OCR.

**Text Referring.** Determining the location of texts accurately is necessary for real-world OCR applications. This ability is evaluated with text grounding and VQA with position tasks.

**Text Spotting.** Text spotting is a widely studied OCR task that requires models to output both the location and content of text. We consider it a distinct capability due to this unique output format.

**Relation Extraction.** Given that texts are often densely arranged in images, the ability to extract and map visual components is essential. This capability is assessed through key information extraction, key information mapping, and handwritten content extraction.

**Element Parsing.** LMMs face the need of parsing complex elements for downstream applications. This ability is evaluated via table parsing, chart parsing, document parsing, and formula recognition.

**Mathematical Calculation.** Math calculation is essential for LMMs to address numerical reasoning tasks. Hence, text counting is introduced to assess the textual perception ability. Besides, we enhance the math QA data by rendering textual questions into images, accompanied by geometric figures.

**Visual Text Understanding.** To tackle sophisticated tasks involving human interaction, LMMs need to comprehend the semantic information of texts, a capability we term visual text understanding. This ability is evaluated by document classification and diagram QA. Additionally, we include basic VQA instructions where answers are located directly within the image, which refers to cognition VQA.

**Knowledge Reasoning.** Some tasks require complex inference and world knowledge, including science QA, APP agent interactions, ASCII art classification, text translation, and reasoning VQA (where answers are not directly visible in images).

## 4.2 Annotation Curation

**Dataset Collection.** To ensure data diversity, we manually harvest and screen 81 text-rich academic datasets. To ensure diverse scenario coverage, we also supplement them with additional private data. In all, our dataset comprises 31 typical scenarios (see Tab. 11 for the full list).

**Annotation Protocol.** Before starting the annotation, we conducted thorough discussions to establish clear guidelines. For example, in questions involving numbers such as dates, amounts, or frequencies, answers were required to include all common formats—Arabic numerals, English abbreviations, and full English expressions. For coordinate-related questions, all coordinate values in the answers were normalized to a 0–1000 scale based on the image size to ensure consistency across varying image resolutions. In cases where multiple correct answers were possible, all valid answers were included. For the "read all text" task, we required that the answer follow a natural reading order from left to right and from top to bottom. Based on these guidelines, 15 professional annotators carried out the annotation work. Each annotator strictly adhered to the instructions and created QA pairs along with the relevant coordinate information, depending on the task requirements.

**Manual Verification.** To ensure data quality, we perform a manual cross-validation process to ensure accuracy and quality. Specifically, each annotated example was first completed by one annotator, then reviewed by a second annotator to verify the correctness. If disagreements or ambiguities arose, the case was escalated to a third annotator for judgment. In instances where consensus could not be reached among all three annotators, the corresponding instruction was excluded from the dataset. Finally approximately 1% annotations are corrected.

## 4.3 Statistics of OCRBench v2

Here we present the OCR-related statistics and the measurement of prompt quality. As shown in Fig. 4 (a) and (b), we count the distribution of line-level OCR results of $7,400$ English and $2,600$ Chinese images. And Fig. 4 (c) exhibits the average number of line-level OCR results per category. These statistics demonstrate that the text information is sufficiently rich in *OCRBench v2*. In addition, Fig. 4 (d) compares the Average Entropy, Type-Token Ratio, and Average Variability Index of the
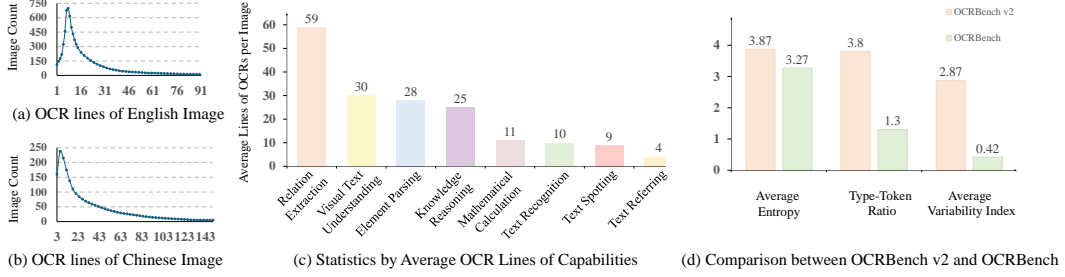
(a) OCR lines of English Image

(b) OCR lines of Chinese Image

(c) Statistics by Average OCR Lines of Capabilities

(d) Comparison between OCRBench v2 and OCRBench

Figure 4: **OCR-related statistics and prompt quality assessment of OCRBench v2**.

Table 2: **Evaluation of existing LMMs on English tasks of OCRBench v2's public data**. "Recognition", "Referring", "Spotting", "Extraction", "Parsing", "Calculation", "Understanding", and "Reasoning" refer to text recognition, text referring, text spotting, relation extraction, element parsing, mathematical calculation, visual text understanding, and knowledge reasoning, respectively. Higher values indicate better performance. Best performance is in boldface, and the second best is underlined. The notations apply to all subsequent figures.

| Method | Recognition | Referring | Spotting | Extraction | Parsing | Calculation | Understanding | Reasoning | Average |
|---|---|---|---|---|---|---|---|---|---|
| Open-source LMMs | | | | | | | | | |
| LLaVA-Next-8B [49] | 41.3 | 18.8 | 0 | 49.5 | 21.2 | 17.3 | 55.2 | 48.9 | 31.5 |
| LLaVA-OV-7B [50] | 46.0 | 20.8 | 0.1 | 58.3 | 25.3 | 23.3 | 64.4 | 53.0 | 36.4 |
| Monkey [51] | 35.2 | 0 | 0 | 16.6 | 16.3 | 14.4 | 59.8 | 42.3 | 23.1 |
| TextMonkey [7] | 39.1 | 0.7 | 0 | 19.0 | 12.2 | 19.0 | 61.1 | 40.2 | 23.9 |
| Molmo-7B [52] | 52.4 | 21.3 | 0.1 | 45.5 | 7.6 | 28.5 | 65.3 | 55.0 | 34.5 |
| Cambrian-1-8B [53] | 45.3 | 21.5 | 0 | 53.6 | 19.2 | 19.5 | 63.5 | 55.5 | 34.7 |
| Pixtral-12B [54] | 48.9 | 21.6 | 0 | 66.3 | 35.5 | 29.8 | 66.9 | 53.7 | 40.3 |
| Qwen2.5-VL-7B [21] | 68.8 | 25.7 | 1.2 | 80.2 | 30.4 | 38.2 | 73.2 | 56.2 | 46.7 |
| InternVL3-14B [23] | 67.3 | 36.9 | 11.2 | 89.0 | 38.4 | 38.4 | 79.2 | 60.5 | 52.6 |
| Deepseek-VL2-Small [55] | 62.7 | 28.0 | 0.1 | 77.5 | 32.7 | 14.3 | 77.1 | 53.9 | 43.3 |
| MiniCPM-o-2.6 [56] | 66.9 | 29.5 | 0.5 | 70.8 | 33.4 | 31.9 | 69.9 | 57.9 | 45.1 |
| GLM-4V-9B [57] | 61.8 | 22.6 | 0 | 71.7 | 31.6 | 22.6 | 72.1 | 58.4 | 42.6 |
| Ovis2-8B [58] | 73.2 | 24.6 | 0.7 | 62.4 | 44.8 | 40.6 | 72.7 | 62.6 | 47.7 |
| Closed-source LMMs | | | | | | | | | |
| GPT-4o [1] | 61.2 | 26.7 | 0 | 77.5 | 36.3 | 43.4 | 71.1 | 55.5 | 46.5 |
| GPT-4o-mini [59] | 57.9 | 23.3 | 0.6 | 70.8 | 31.5 | 38.8 | 65.9 | 55.1 | 43.0 |
| Gemini-Pro [60] | 61.2 | 39.5 | 13.5 | 79.3 | 39.2 | 47.7 | 75.5 | 59.3 | 51.9 |
| Claude3.5-sonnet [61] | 62.2 | 28.4 | 1.3 | 56.6 | 37.8 | 40.8 | 73.5 | 60.9 | 45.2 |
| Step-1V [62] | 67.8 | 31.3 | 7.2 | 73.6 | 37.2 | 27.8 | 69.8 | 58.6 | 46.7 |

questions between *OCRBench v2* and OCRBench. *OCRBench v2* presents higher values across all three metrics, indicating more diverse, less redundant, and structurally varied questions. This suggests it provides a more comprehensive and challenging benchmark for LMMs.

## 4.4 Evaluation Criteria

We adopt six types of evaluation metrics tailored to specific task categories. In the following, we present an overview of the evaluation metrics and their applicability to specific tasks.

**Parsing Type.** To evaluate the element parsing ability of LMMs, we assess their performance in transforming input images into structured formats, including HTML, Markdown, and JSON. TEDS [63] is employed to measure the structural similarity between outputs and the desired format.

**Localization Type.** For text referring, the IoU score is applied to quantify the distance between the predicted regions and the ground truth.

**Extraction Type.** To evaluate relation extraction, we employ the F1 score to assess key information extraction and mapping. Since this evaluation requires structural extraction of information from the output of LMMs, the format is provided in the given prompt.

**Long Reading Type.** To assess performance on long text reading tasks, BLEU [64], METEOR [65], F1 score, and edit distance are used to assess the similarity between predicted text and ground truth.

Table 3: **Evaluation of existing LMMs on Chinese tasks of OCRBench v2's public data**. "LLM Size" indicates the number of parameters of the language model employed in each method.

| Method | LLM Size | Recognition | Extraction | Parsing | Understanding | Reasoning | Average |
|---|---|---|---|---|---|---|---|
| Open-source LMMs | | | | | | | |
| LLaVA-Next-8B [49] | 8B | 5.7 | 2.9 | 12.2 | 7.5 | 17.2 | 9.1 |
| LLaVA-OV-7B [50] | 8B | 14.8 | 15.7 | 13.7 | 16.0 | 28.7 | 17.8 |
| Monkey [51] | 8B | 4.6 | 11.2 | 8.4 | 21.5 | 20.0 | 13.1 |
| TextMonkey [7] | 8B | 23.5 | 14.8 | 8.4 | 19.9 | 12.2 | 15.8 |
| Molmo-7B [52] | 8B | 7.1 | 15.0 | 9.2 | 9.0 | 23.7 | 12.8 |
| Cambrian-1-8B [53] | 8B | 5.3 | 14.9 | 12.6 | 8.5 | 8.1 | 9.9 |
| Pixtral-12B [54] | 12B | 13.4 | 10.9 | 21.0 | 7.0 | 20.7 | 14.6 |
| Qwen2.5-VL-7B [21] | 8B | **75.3** | 61.4 | **41.8** | 59.3 | 40.4 | 55.6 |
| InternVL3-14B [23] | 14B | 66.2 | **64.8** | 33.5 | **63.4** | **50.6** | **55.7** |
| Deepseek-VL2-Small [55] | 16B | 60.9 | 50.6 | 28.3 | 53.0 | 20.5 | 42.7 |
| MiniCPM-o-2.6 [56] | 7B | 53.0 | 49.4 | 27.1 | 43.5 | 32.7 | 41.1 |
| GLM-4V-9B [57] | 9B | 24.4 | 60.6 | 20.4 | 52.8 | 25.2 | 36.6 |
| Ovis2-8B [58] | 7B | 72.2 | 50.8 | 37.7 | 47.9 | 37.4 | 49.2 |
| Closed-source LMMs | | | | | | | |
| GPT-4o [1] | - | 21.6 | 53.0 | 29.8 | 38.5 | 18.2 | 32.2 |
| GPT-4o-mini [59] | - | 13.1 | 38.9 | 27.2 | 28.8 | 16.9 | 25.0 |
| Gemini-Pro [60] | - | 52.5 | 47.3 | 30.9 | 51.5 | 33.4 | 43.1 |
| Claude3.5-sonnet [61] | - | 21.0 | 56.2 | 35.2 | 55.0 | 30.5 | 39.6 |
| Step-1V [62] | - | 56.7 | 41.1 | 37.6 | 38.3 | 39.2 | 42.6 |

**Counting Type.** In text counting, LMMs are required to count the number of text instances. Thus, we use the L1 distance to measure the absolute difference between predicted and ground truth counts. The final score is then normalized to the range of $[0, 1]$ based on the ground truth.

**Basic VQA Type.** For questions where the original data provides options, we use exact string matching to compute accuracy. In other cases, we follow the approach of OCRBench to check whether the ground truth is contained in the prediction for short answers (fewer than 5 words) and employ ANLS to measure prediction quality for longer answers (5 words or more).

# 5 Results and Findings

Here we first benchmark state-of-the-art LMMs on *OCRBench v2*, presenting the quantitative analysis, then summarize key findings of current limitations for LMMs. All results are presented as percentages.

## 5.1 Baselines

The tested LMMs in this section includes LLaVA-Next-8B [49], LLaVA-OV-7B [50], Monkey [51], TextMonkey [7], Molmo-7B [52], Cambrian-1-8B [53], Pixtral-12B [54], Qwen2.5-VL-7B [21], InternVL3-14B [23], Deepseek-VL2-Tiny [55], MiniCPM-o-2.6 [56], GLM-4v-9B [57], Ovis2-8B [58], GPT4o [24], GPT4o-mini [59], Gemini-1.5-Pro [60], Claude3.5-sonnet [61], and Step-1V [62]. More LMM evaluation results can be found in Tabs. 12, 13, 14, and 15.

## 5.2 Main Results

**Evaluation results on public data** are shown in Tab. 2 and Tab. 3. While LMMs perform well on some basic capabilities such as text recognition and visual text understanding, most LMMs achieve low scores in other capabilities, such as text spotting and element parsing, mostly below 50. In particular, some LMMs show significant limitations in text spotting capabilities, failing to precisely locate and recognize the texts. Additionally, LMMs demonstrate inadequate abilities in element parsing and mathematical calculation, which are crucial for complicated tasks like document analysis and mathematical reasoning. Besides, after comparing the performance of LMMs on visual text understanding and knowledge reasoning capabilities, we find that they perform poorly in knowledge reasoning. This suggests the deficiency of LMMs in logical reasoning.

**Evaluation results on private data** are shown in Tab. 4 and Tab. 5. We observe similar evaluation trends to those in the public test set experiments. Overall, LMMs exhibit unsatisfactory performance in text referring, text spotting, element parsing, mathematical calculation, and knowledge reasoning capabilities. In addition, closed-source LMMs outperform their open-source counterparts, demon-

Table 4: **Evaluation of existing LMMs on English tasks of OCRBench v2's private data**.

| Method | Recognition | Referring | Spotting | Extraction | Parsing | Calculation | Understanding | Reasoning | Average |
|---|---|---|---|---|---|---|---|---|---|
| Open-source LMMs | | | | | | | | | |
| LLaVA-Next-8B [49] | 41.4 | 17.0 | 0 | 49.0 | 12.9 | 16.1 | 60.9 | 30.5 | 28.5 |
| LLaVA-OV-7B [50] | 45.4 | 18.5 | 0 | 60.0 | 15.5 | 32.0 | 59.0 | 39.3 | 33.7 |
| Monkey [51] | 31.5 | 0.1 | 0 | 34.4 | _26.3_ | 17.7 | 61.4 | 22.4 | 24.2 |
| TextMonkey [7] | 39.8 | 1.6 | 0 | 27.6 | 24.8 | 10.2 | 62.3 | 21.2 | 23.4 |
| Molmo-7B [52] | 40.8 | 19.5 | 0 | 51.7 | 10.0 | 33.9 | 67.0 | 48.0 | 33.9 |
| Cambrian-1-8B [53] | 44.0 | 19.0 | 0 | 52.3 | 19.0 | 20.7 | 64.0 | 39.3 | 32.3 |
| Pixtral-12B [54] | 45.1 | 21.8 | 0 | 71.6 | 21.7 | 30.4 | 77.3 | 39.5 | 38.4 |
| Qwen2.5-VL-7B [66] | 51.5 | 24.5 | _3.1_ | 64.8 | 13.1 | 53.3 | _78.6_ | 45.5 | 41.8 |
| InternVL3-14B [23] | 55.8 | 24.5 | 2.1 | _89.3_ | 21.0 | _59.5_ | 72.0 | 50.0 | 46.8 |
| Deepseek-VL2-Small [55] | 56.6 | 23.7 | 0 | 86.4 | 18.9 | 30.6 | 72.2 | 39.5 | 41.0 |
| MiniCPM-o-2.6 [56] | 54.1 | 24.7 | 0.3 | 74.4 | 17.6 | 39.2 | 75.7 | 47.0 | 41.6 |
| GLM-4v-9B [57] | 52.7 | 20.6 | 0 | 79.4 | 15.9 | 21.5 | 74.7 | 32.0 | 37.1 |
| Ovis2-8B [58] | 54.2 | 20.9 | 0 | 83.6 | 24.2 | 54.7 | 74.1 | 57.3 | 46.1 |
| Closed-source LMMs | | | | | | | | | |
| GPT-4o [1] | _58.6_ | 23.4 | 0 | 87.4 | 23.1 | 51.6 | 74.4 | **62.3** | _47.6_ |
| GPT-4o-mini [59] | 55.3 | 21.8 | 0 | 85.4 | 20.6 | 45.2 | 75.5 | 49.0 | 44.1 |
| Gemini1.5-Pro [60] | **59.1** | **41.2** | **6.6** | **89.5** | 22.4 | 54.7 | **78.8** | _60.3_ | **51.6** |
| Claude3.5-sonnet [61] | 52.9 | 24.9 | 2.5 | 86.9 | 23.8 | **61.4** | 74.4 | 53.0 | 47.5 |
| Step-1V [62] | 56.7 | _27.4_ | 2.6 | 86.3 | **33.3** | 42.6 | 76.6 | 48.7 | 46.8 |

Table 5: **Evaluation of existing LMMs on Chinese tasks of OCRBench v2's private data**.

| Method | LLM Size | Recognition | Extraction | Parsing | Understanding | Reasoning | Average |
|---|---|---|---|---|---|---|---|
| Open-source LMMs | | | | | | | |
| LLaVA-Next-8B [49] | 8B | 2.8 | 0.9 | 14.9 | 20.0 | 7.4 | 9.2 |
| LLaVA-OV-7B [50] | 8B | 5.4 | 13.6 | 20.3 | 34.0 | 13.6 | 17.4 |
| Monkey [51] | 8B | 1.5 | 28.4 | 29.1 | 40.0 | 8.3 | 21.5 |
| TextMonkey [7] | 8B | 10.5 | 15.2 | 30.2 | 44.0 | 7.6 | 21.5 |
| Molmo-7B [52] | 8B | 3.4 | 29.8 | 6.6 | 24.0 | 11.1 | 15.0 |
| Cambrian-1-8B [53] | 8B | 2.4 | 19.8 | 26.7 | 36.0 | 7.6 | 18.5 |
| Pixtral-12B [54] | 12B | 6.2 | 22.3 | 11.4 | 26.0 | 14.0 | 16.0 |
| Qwen2.5-VL-7B [66] | 8B | 24.4 | **78.9** | 33.1 | **82.0** | 29.0 | 49.5 |
| InternVL3-14B [23] | 14B | 62.1 | 59.5 | 33.2 | 80.0 | 29.2 | 52.8 |
| DeepSeek-VL2-Small [55] | 16B | 51.6 | 56.3 | 27.8 | 79.6 | 25.3 | 48.1 |
| MiniCPM-o-2.6 [56] | 7B | 54.0 | 62.4 | 24.1 | 68.0 | 29.8 | 47.7 |
| GLM-4v-9B [57] | 9B | 60.6 | 65.2 | 32.4 | **82.0** | 18.2 | 51.7 |
| Ovis2-8B [58] | 7B | 61.0 | _67.7_ | **43.6** | **82.0** | 25.6 | **56.0** |
| Closed-source LMMs | | | | | | | |
| GPT-4o [1] | - | 41.7 | 52.1 | 29.0 | 76.0 | 29.4 | 45.7 |
| GPT-4o-mini [59] | - | 20.0 | 53.6 | 27.9 | 66.0 | 19.6 | 37.4 |
| Gemini1.5-Pro [60] | - | **71.4** | 63.8 | 30.5 | **82.0** | _29.9_ | _55.5_ |
| Claude3.5-sonnet [61] | - | 34.2 | 62.5 | _35.2_ | 78.0 | **32.2** | 48.4 |
| Step-1V [62] | - | _65.2_ | 64.9 | 33.1 | 78.0 | 25.5 | 53.4 |

strating stronger generalization capabilities. The consistent results across both public and private test sets confirm the soundness of *OCRBench v2*'s task design, data collection process, and evaluation metrics, and demonstrate its effectiveness in revealing the capability limitations of current LMMs.

## 5.3 Main Findings

We provide in-depth analyses for LMMs' common limitations, including rare text recognition, fine-grained spatial perception, layout perception, complex element analysis, and logical reasoning.

**Finding 1.** LMMs still face challenges with less frequently encountered texts, such as dot matrix texts and mathematical formulas. This performance gap highlights the continuing challenges LMMs face in real-world text recognition. For instance, occluded text, CAPTCHA, and dot-matrix text are considered low-frequency text, whereas other types belong to high-frequency text. Tab. 6 shows the performance of some LMMs on high-frequency and low-frequency texts. Notably, recognition accuracy varies significantly across these categories. For example, InternVL3-14B achieves 79.1% on high-frequency texts but drops to 46.7% on low-frequency ones.

**Finding 2.** Current LMMs still exhibit limited performance in tasks requiring precise spatial understanding, such as text referring and text spotting. For instance, when provided with coordinate information as input, many models are able to output the relevant content from captions or chapters.

Table 6: **LMMs' performance on high- and low-frequency words**.

| Category | Pixtral-12B [54] | Cambrian-1-8B [53] | InternVL3-14B [23] | Qwen2.5-VL-7B [66] |
|---|---|---|---|---|
| High Frequency | 58.3 | 59.8 | 79.1 | 84.5 |
| Low Frequency | 23.6 | 40.2 | 46.7 | 53.3 |

However, almost all models struggle to accurately retrieve the corresponding text from documents with dense text based on given coordinates. We investigate the content response accuracy and the IoU score for answer region localization in the VQA with position task. Tab. 7 suggests that although LMMs can roughly identify where the answer is located, they struggle to output the exact region.

**Finding 3.** While LMMs achieve good performance on basic text recognition, they struggle with complex layouts such as overlapping or rotated texts. For example, GPT-4o fails to detect the characters in overlapping handwritten text and misrecognizes numbers in 90° rotated images, revealing LMMs' limitations in handling texts with complex layouts. Rotating images in the DocVQA dataset led to a significant performance drop of $55.7\%$ for InternVL3-14B (from $90.9\%$ to $35.2\%$).

Table 7: **LMMs' performance on VQA with position task**.

| Category | Pixtral-12B [54] | Cambrian-1-8B [53] | InternVL3-14B [23] | Qwen2.5-VL-7B [66] |
|---|---|---|---|---|
| Content Accuarcy | 68.8 | 71.7 | 78.3 | 75.2 |
| IoU Accuracy | 1.7 | 0.0 | 12.9 | 9.6 |

**Finding 4.** LMMs still struggle to parse text into structured formats in downstream applications such as document digitalization. For instance, InternVL3-14B achieves 94.4% accuracy in unpaired entities matching, but its performance drops to 84.9% in key information extraction, where the model is required to identify the corresponding value given an entity. The performance further degrades in element parsing tasks that demand structured outputs.

**Finding 5.** Despite recent advances, LMMs still face challenges in complex mathematical and textual reasoning tasks. To assess their capabilities, we evaluated InternVL3-14B on the private test set covering reasoning VQA, ScienceQA, and APP agent tasks. Questions were categorized into five types: common sense reasoning, visual-text understanding, pattern recognition, calculation, and expert knowledge. Human ratings showed the model achieved accuracies of 72.9%, 83.0%, 69.2%, 56.5%, and 71.8%, respectively, indicating notable variation.

# 6   Conclusion

In this work, we introduce *OCRBench v2*, a comprehensive benchmark designed to evaluate the OCR capabilities of LMMs. Covering 23 tasks across 31 diverse scenarios, our benchmark systematically assesses eight core capabilities that are essential for text-oriented visual understanding tasks. It includes $10,000$ high-quality QA pairs and six rigorous evaluation metrics. In addition, we curate a private test set of $1,500$ manually labeled images to ensure robust generalization evaluation. Leveraging this benchmark, we conduct extensive experiments on representative LMMs. Through in-depth analysis of experimental results, we identify critical limitations of current models and uncover key factors that affect their OCR performance. We hope *OCRBench v2* could aid future research on enhancing LMMs' text understanding ability.

# Acknowledgements

# References

[1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.

[2] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.

[3] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in Neural Information Processing Systems*, 2020.

[4] J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, and J. Zhou, "Qwen-vl: A frontier large vision-language model with versatile abilities," *arXiv preprint arXiv:2308.12966*, 2023.

[5] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[6] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny, "Minigpt-4: Enhancing vision-language understanding with advanced large language models," *Proceedings of the International Conference on Learning Representations*, 2024.

[7] Y. Liu, B. Yang, Q. Liu, Z. Li, Z. Ma, S. Zhang, and X. Bai, "Textmonkey: An ocr-free large multimodal model for understanding document," *arXiv preprint arXiv:2403.04473*, 2024.

[8] C. Fu, P. Chen, Y. Shen, Y. Qin, M. Zhang, X. Lin, J. Yang, X. Zheng, K. Li, X. Sun *et al.*, "MME: A comprehensive evaluation benchmark for multimodal large language models," *arXiv preprint arXiv:2306.13394*, 2023.

[9] K. Ying, F. Meng, J. Wang, Z. Li, H. Lin, Y. Yang, H. Zhang, W. Zhang, Y. Lin, S. Liu *et al.*, "Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi," *arXiv preprint arXiv:2404.16006*, 2024.

[10] A. Singh, V. Natarajan, M. Shah, Y. Jiang, X. Chen, D. Batra, D. Parikh, and M. Rohrbach, "Towards vqa models that can read," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8317–8326.

[11] A. F. Biten, R. Tito, A. Mafla, L. Gomez, M. Rusinol, E. Valveny, C. Jawahar, and D. Karatzas, "Scene text visual question answering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4291–4301.

[12] X. Wang, Y. Liu, C. Shen, C. C. Ng, C. Luo, L. Jin, C. S. Chan, A. v. d. Hengel, and L. Wang, "On the general value of evidence, and bilingual scene-text visual question answering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 126–10 135.

[13] L. Chen, J. Li, X. Dong, P. Zhang, Y. Zang, Z. Chen, H. Duan, J. Wang, Y. Qiao, D. Lin *et al.*, "Are We on the Right Way for Evaluating Large Vision-Language Models?" *arXiv preprint arXiv:2403.20330*, 2024.

[14] Y. Liu, Z. Li, B. Yang, C. Li, X. Yin, C.-l. Liu, L. Jin, and X. Bai, "On the hidden mystery of ocr in large multimodal models," *arXiv preprint arXiv:2305.07895*, 2023.

[15] B. Li, Y. Ge, Y. Chen, Y. Ge, R. Zhang, and Y. Shan, "Seed-bench-2-plus: Benchmarking multimodal large language models with text-rich visual comprehension," *arXiv preprint arXiv:2404.16790*, 2024.

[16] R. Wadhawan, H. Bansal, K.-W. Chang, and N. Peng, "ConTextual: Evaluating Context-Sensitive Text-Rich Visual Reasoning in Large Multimodal Models," in *Proceedings of International Conference on Machine Learning*, 2024.

[17] C. Liu, H. Wei, J. Chen, L. Kong, Z. Ge, Z. Zhu, L. Zhao, J. Sun, C. Han, and X. Zhang, "Focus Anywhere for Fine-grained Multi-page Document Understanding," *arXiv preprint arXiv:2405.14295*, 2024.

[18] Y. Kim, M. Yim, and K. Y. Song, "TableVQA-Bench: A visual question answering benchmark on multiple table domains," *arXiv preprint arXiv:2404.19205*, 2024.

[19] W. Zhao, H. Feng, Q. Liu, J. Tang, S. Wei, B. Wu, L. Liao, Y. Ye, H. Liu, H. Li *et al.*, "TabPedia: Towards Comprehensive Visual Table Understanding with Concept Synergy," *arXiv preprint arXiv:2406.01326*, 2024.

[20] R. Xia, B. Zhang, H. Ye, X. Yan, Q. Liu, H. Zhou, Z. Chen, M. Dou, B. Shi, J. Yan *et al.*, "Chartx & chartvlm: A versatile benchmark and foundation model for complicated chart reasoning," *arXiv preprint arXiv:2402.12185*, 2024.

[21] Q. Team, "Qwen2.5-vl," January 2025. [Online]. Available: https://qwenlm.github.io/blog/qwen2.5-vl/

[22] M. Mathew, D. Karatzas, and C. Jawahar, "Docvqa: A dataset for vqa on document images," in *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*, 2021, pp. 2200–2209.

[23] Z. Chen, W. Wang, Y. Cao, Y. Liu, Z. Gao, E. Cui, J. Zhu, S. Ye, H. Tian, Z. Liu *et al.*, "Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling," *arXiv preprint arXiv:2412.05271*, 2024.

[24] OpenAI, "Hello GPT-4o," https://openai.com/index/gpt-4v-system-card, 2024, accessed: 2024-12-29.

[25] L. Ouyang, Y. Qu, H. Zhou, J. Zhu, R. Zhang, Q. Lin, B. Wang, Z. Zhao, M. Jiang, X. Zhao *et al.*, "Omnidocbench: Benchmarking diverse pdf document parsing with comprehensive annotations," *arXiv preprint arXiv:2412.07626*, 2024.

[26] Z. Yang, J. Tang, Z. Li, P. Wang, J. Wan, H. Zhong, X. Liu, M. Yang, P. Wang, Y. Liu *et al.*, "Cc-ocr: A comprehensive and challenging ocr benchmark for evaluating large multimodal models in literacy," *arXiv preprint arXiv:2412.02210*, 2024.

[27] Y. Ma, Y. Zang, L. Chen, M. Chen, Y. Jiao, X. Li, X. Lu, Z. Liu, Y. Ma, X. Dong *et al.*, "Mmlongbench-doc: Benchmarking long-context document understanding with visualizations," *arXiv preprint arXiv:2407.01523*, 2024.

[28] M. Zheng, X. Feng, Q. Si, Q. She, Z. Lin, W. Jiang, and W. Wang, "Multimodal Table Understanding," in *Proceedings of Annual Meeting of the Association for Computational Linguistics*, L. Ku, A. Martins, and V. Srikumar, Eds. Association for Computational Linguistics, 2024, pp. 9102–9124. [Online]. Available: https://doi.org/10.18653/v1/2024.acl-long.493

[29] F. Liu, X. Wang, W. Yao, J. Chen, K. Song, S. Cho, Y. Yacoob, and D. Yu, "MMC: Advancing Multimodal Chart Understanding with Large-scale Instruction Tuning," in *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2024, pp. 1287–1310.

[30] Y. Zhang, R. Zhang, J. Gu, Y. Zhou, N. Lipka, D. Yang, and T. Sun, "Llavar: Enhanced visual instruction tuning for text-rich image understanding," *arXiv preprint arXiv:2306.17107*, 2023.

[31] J. Ye, A. Hu, H. Xu, Q. Ye, M. Yan, Y. Dan, C. Zhao, G. Xu, C. Li, J. Tian *et al.*, "mplug-docowl: Modularized multimodal large language model for document understanding," *arXiv preprint arXiv:2307.02499*, 2023.

[32] H. Feng, Q. Liu, H. Liu, W. Zhou, H. Li, and C. Huang, "Docpedia: Unleashing the power of large multimodal model in the frequency domain for versatile document understanding," *arXiv preprint arXiv:2311.11810*, 2023.

[33] J. Ye, A. Hu, H. Xu, Q. Ye, M. Yan, G. Xu, C. Li, J. Tian, Q. Qian, J. Zhang *et al.*, "Ureader: Universal ocr-free visually-situated language understanding with multimodal large language model," *arXiv preprint arXiv:2310.05126*, 2023.

[34] C. Luo, Y. Shen, Z. Zhu, Q. Zheng, Z. Yu, and C. Yao, "LayoutLLM: Layout Instruction Tuning with Large Language Models for Document Understanding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 15 630–15 640.

[35] A. Hu, H. Xu, J. Ye, M. Yan, L. Zhang, B. Zhang, C. Li, J. Zhang, Q. Jin, F. Huang *et al.*, "mplug-docowl 1.5: Unified structure learning for ocr-free document understanding," *arXiv preprint arXiv:2403.12895*, 2024.

[36] J. Zhang, W. Yang, S. Lai, Z. Xie, and L. Jin, "Dockylin: A large multimodal model for visual document understanding with efficient visual slimming," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 39, no. 9, 2025, pp. 9923–9932.

[37] W. Liao, J. Wang, H. Li, C. Wang, J. Huang, and L. Jin, "Doclayllm: An efficient and effective multi-modal extension of large language models for text-rich document understanding," *arXiv preprint arXiv:2408.15045*, 2024.

[38] Z. Zhu, C. Luo, Z. Shao, F. Gao, H. Xing, Q. Zheng, and J. Zhang, "A simple yet effective layout token in large language models for document understanding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025.

[39] H. Xiao, Y. Xie, G. Tan, Y. Chen, R. Hu, K. Wang, A. Zhou, H. Li, H. Shao, X. Lu, P. Gao, Y. Wen, X. Chen, S. Ren, and H. Li, "Adaptive markup language generation for contextually-grounded visual document understanding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025.

[40] Z. Wang, T. Guan, P. Fu, C. Duan, Q. Jiang, Z. Guo, S. Guo, J. Luo, W. Shen, and X. Yang, "Marten: Visual question answering with mask generation for multi-modal document understanding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025.

[41] A. Masry, D. X. Long, J. Q. Tan, S. Joty, and E. Hoque, "Chartqa: A benchmark for question answering about charts with visual and logical reasoning," *arXiv preprint arXiv:2203.10244*, 2022.

[42] M. Mathew, V. Bagal, R. Tito, D. Karatzas, E. Valveny, and C. Jawahar, "Infographicvqa," in *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*, 2022, pp. 1697–1706.

[43] S. Zhang, B. Yang, Z. Li, Z. Ma, Y. Liu, and X. Bai, "Exploring the Capabilities of Large Multimodal Models on Dense Text," in *Proceedings of International Conference on Document Analysis and Recognition*. Springer, 2024, pp. 281–298.

[44] J. Chen, L. Kong, H. Wei, C. Liu, Z. Ge, L. Zhao, J. Sun, C. Han, and X. Zhang, "Onechart: Purify the chart structural extraction via one auxiliary token," in *Proceedings of the ACM International Conference on Multimedia*, 2024, pp. 147–155.

[45] J. Van Landeghem, R. Tito, Ł. Borchmann, M. Pietruszka, P. Joziak, R. Powalski, D. Jurkiewicz, M. Coustaty, B. Anckaert, E. Valveny *et al.*, "Document understanding dataset and evaluation (dude)," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 19 528–19 540.

[46] W. Wang, S. Zhang, Y. Ren, Y. Duan, T. Li, S. Liu, M. Hu, Z. Chen, K. Zhang, L. Lu *et al.*, "Needle in a multimodal haystack," *Advances in Neural Information Processing Systems*, vol. 37, pp. 20 540–20 565, 2025.

[47] R. Tito, D. Karatzas, and E. Valveny, "Hierarchical multimodal transformers for multipage docvqa," *Pattern Recognition*, vol. 144, p. 109834, 2023.

[48] C. Deng, J. Yuan, P. Bu, P. Wang, Z.-Z. Li, J. Xu, X.-H. Li, Y. Gao, J. Song, B. Zheng *et al.*, "Longdocurl: a comprehensive multimodal long document benchmark integrating understanding, reasoning, and locating," *arXiv preprint arXiv:2412.18424*, 2024.

[49] H. Liu, C. Li, Y. Li, B. Li, Y. Zhang, S. Shen, and Y. J. Lee, "Llava-next: Improved reasoning, ocr, and world knowledge," 2024.

[50] B. Li, Y. Zhang, D. Guo, R. Zhang, F. Li, H. Zhang, K. Zhang, Y. Li, Z. Liu, and C. Li, "Llava-onevision: Easy visual task transfer," *arXiv preprint arXiv:2408.03326*, 2024.

[51] Z. Li, B. Yang, Q. Liu, Z. Ma, S. Zhang, J. Yang, Y. Sun, Y. Liu, and X. Bai, "Monkey: Image resolution and text label are important things for large multi-modal models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 26 763–26 773.

[52] M. Deitke, C. Clark, S. Lee, R. Tripathi, Y. Yang, J. S. Park, M. Salehi, N. Muennighoff, K. Lo, L. Soldaini *et al.*, "Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models," *arXiv preprint arXiv:2409.17146*, 2024.

[53] S. Tong, E. L. Brown II, P. Wu, S. Woo, A. J. IYER, S. C. Akula, S. Yang, J. Yang, M. Middepogu, Z. Wang *et al.*, "Cambrian-1: A fully open, vision-centric exploration of multimodal llms," in *Advances in Neural Information Processing Systems*, 2024.

[54] P. Agrawal, S. Antoniak, E. B. Hanna, B. Bout, D. Chaplot, J. Chudnovsky, D. Costa, B. De Monicault, S. Garg, T. Gervet *et al.*, "Pixtral 12b," *arXiv preprint arXiv:2410.07073*, 2024.

[55] Z. Wu, X. Chen, Z. Pan, X. Liu, W. Liu, D. Dai, H. Gao, Y. Ma, C. Wu, B. Wang *et al.*, "Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding," *arXiv preprint arXiv:2412.10302*, 2024.

[56] Y. Yao, T. Yu, A. Zhang, C. Wang, J. Cui, H. Zhu, T. Cai, H. Li, W. Zhao, Z. He *et al.*, "Minicpm-v: A gpt-4v level mllm on your phone," *arXiv preprint arXiv:2408.01800*, 2024.

[57] T. GLM, A. Zeng, B. Xu, B. Wang, C. Zhang, D. Yin, D. Rojas, G. Feng, H. Zhao, H. Lai *et al.*, "ChatGLM: A Family of Large Language Models from GLM-130B to GLM-4 All Tools," *arXiv preprint arXiv:2406.12793*, 2024.

[58] S. Lu, Y. Li, Q.-G. Chen, Z. Xu, W. Luo, K. Zhang, and H.-J. Ye, "Ovis: Structural embedding alignment for multimodal large language model," *arXiv preprint arXiv:2405.20797*, 2024.

[59] OpenAI, "GPT-4o mini: advancing cost-efficient intelligence," https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence, 2024, accessed: 2024-12-29.

[60] G. Team, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, K. Millican *et al.*, "Gemini: a family of highly capable multimodal models," *arXiv preprint arXiv:2312.11805*, 2023.

[61] Anthropic, "Claude 3.5 Sonnet," https://www.anthropic.com/news/claude-3-5-sonnet, 2024, accessed: 2024-12-29.

[62] StepFun, "Step-1V," https://www.stepfun.com/#step1v, 2024, accessed: 2024-12-29.

[63] X. Zhong, E. ShafieiBavani, and A. Jimeno Yepes, "Image-based table recognition: data, model, and evaluation," in *Proceedings of European Conference on Computer Vision*. Springer, 2020, pp. 564–580.

[64] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of Annual Meeting of the Association for Computational Linguistics*, 2002, pp. 311–318.

[65] S. Banerjee and A. Lavie, "METEOR: An automatic metric for mt evaluation with improved correlation with human judgments," in *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, 2005, pp. 65–72.

[66] P. Wang, S. Bai, S. Tan, S. Wang, Z. Fan, J. Bai, K. Chen, X. Liu, J. Wang, W. Ge *et al.*, "Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution," *arXiv preprint arXiv:2409.12191*, 2024.

[67] B. Shi, X. Bai, and C. Yao, "An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 11, pp. 2298–2304, 2016.

[68] S. Fang, H. Xie, Y. Wang, Z. Mao, and Y. Zhang, "Read like humans: Autonomous, bidirectional and iterative language modeling for scene text recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7098–7107.

[69] B. Shi, M. Yang, X. Wang, P. Lyu, C. Yao, and X. Bai, "Aster: An attentional scene text recognizer with flexible rectification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 9, pp. 2035–2048, 2018.

[70] N. Lu, W. Yu, X. Qi, Y. Chen, P. Gong, R. Xiao, and X. Bai, "Master: Multi-aspect non-local network for scene text recognition," *Pattern Recognition*, vol. 117, p. 107980, 2021.

[71] Y. Du, Z. Chen, C. Jia, X. Yin, T. Zheng, C. Li, Y. Du, and Y. Jiang, "SVTR: scene text recognition with a single visual model," in *Proceedings of the International Joint Conference on Artificial Intelligence*, L. D. Raedt, Ed. ijcai.org, 2022, pp. 884–890. [Online]. Available: https://doi.org/10.24963/ijcai.2022/124

[72] Y. Liu, H. Chen, C. Shen, T. He, L. Jin, and L. Wang, "Abcnet: Real-time scene text spotting with adaptive bezier-curve network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9809–9818.

[73] Y. Liu, C. Shen, L. Jin, T. He, P. Chen, C. Liu, and H. Chen, "Abcnet v2: Adaptive bezier-curve network for real-time end-to-end text spotting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 11, pp. 8048–8064, 2021.

[74] X. Zhang, Y. Su, S. Tripathi, and Z. Tu, "Text spotting transformers," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 9519–9528.

[75] C. K. Ch'ng and C. S. Chan, "Total-text: A comprehensive dataset for scene text detection and recognition," in *Proceedings of International Conference on Document Analysis and Recognition*, vol. 1. IEEE, 2017, pp. 935–942.

[76] H. Wei, C. Liu, J. Chen, J. Wang, L. Kong, Y. Xu, Z. Ge, L. Zhao, J. Sun, Y. Peng *et al.*, "General ocr theory: Towards ocr-2.0 via a unified end-to-end model," *arXiv preprint arXiv:2409.01704*, 2024.

[77] D. Karatzas, F. Shafait, S. Uchida, M. Iwamura, L. G. i. Bigorda, S. R. Mestre, J. Mas, D. F. Mota, J. A. Almazàn, and L. P. de las Heras, "Icdar 2013 robust reading competition," in *Proceedings of International Conference on Document Analysis and Recognition*, 2013, pp. 1484–1493.

[78] C. Shi, C. Wang, B. Xiao, S. Gao, and J. Hu, "End-to-end scene text recognition using tree-structured models," *Pattern Recognition*, vol. 47, pp. 2853–2866, 2014. [Online]. Available: https://api.semanticscholar.org/CorpusID:30201169

[79] A. Mishra, K. Alahari, and C. V. Jawahar, "Scene text recognition using higher order language priors," in *British Machine Vision Conference*, 2012.

[80] D. Karatzas, L. Gomez-Bigorda, A. Nicolaou, S. Ghosh, A. Bagdanov, M. Iwamura, J. Matas, L. Neumann, V. R. Chandrasekhar, S. Lu *et al.*, "Icdar 2015 competition on robust reading," in *Proceedings of International Conference on Document Analysis and Recognition*. IEEE, 2015, pp. 1156–1160.

[81] Y. Liu, L. Jin, S. Zhang, C. Luo, and S. Zhang, "Curved scene text detection via transverse and longitudinal sequence connection," *Pattern Recognition*, vol. 90, no. C, p. 337–345, Jun. 2019. [Online]. Available: https://doi.org/10.1016/j.patcog.2019.02.002

[82] A. Veit, T. Matera, L. Neumann, J. Matas, and S. Belongie, "Coco-text: Dataset and benchmark for text detection and recognition in natural images," *arXiv preprint arXiv:1601.07140*, 2016.

[83] A. Risnumawan, P. Shivakumara, C. S. Chan, and C. L. Tan, "A robust arbitrary text detection system for natural scene images," *Expert Systems with Applications*, vol. 41, pp. 8027–8048, 2014. [Online]. Available: https://api.semanticscholar.org/CorpusID:15559857

[84] T. Q. Phan, P. Shivakumara, S. Tian, and C. L. Tan, "Recognizing text with perspective distortion in natural scenes," in *Proceedings of IEEE/CVF International Conference on Computer Vision*. IEEE Computer Society, 2013, pp. 569–576. [Online]. Available: https://doi.org/10.1109/ICCV.2013.76

[85] X. Xie, L. Fu, Z. Zhang, Z. Wang, and X. Bai, "Toward understanding wordart: Corner-guided transformer for scene text recognition," in *Proceedings of European Conference on Computer Vision*, S. Avidan, G. Brostow, M. Cissé, G. M. Farinella, and T. Hassner, Eds. Cham: Springer Nature Switzerland, 2022, pp. 303–321.

[86] U.-V. Marti and H. Bunke, "The iam-database: an english sentence database for offline handwriting recognition," *International Journal on Document Analysis and Recognition*, vol. 5, pp. 39–46, 2002.

[87] M. Diem, S. Fiel, F. Kleber, R. Sablatnig, J. M. Saavedra, D. Contreras, J. M. Barrios, and L. S. Oliveira, "Proceedings of ieee international conference on frontiers in handwriting recognition," in *2014 14th International Conference on Frontiers in Handwriting Recognition*, 2014, pp. 779–784.

[88] Y. Wang, H. Xie, S. Fang, J. Wang, S. Zhu, and Y. Zhang, "From two to one: A new scene text recognizer with visual language modeling network," in *Proceedings of IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14 194–14 203.

[89] L. Yuliang, J. Lianwen, Z. Shuaitao, and Z. Sheng, "Detecting curve text in the wild: New dataset and new solution," *arXiv preprint arXiv:1712.02170*, 2017.

[90] S. Long, S. Qin, D. Panteleev, A. Bissacco, Y. Fujii, and M. Raptis, "Towards end-to-end unified scene text detection and layout analysis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 1049–1059.

[91] T.-L. Yuan, Z. Zhu, K. Xu, C.-J. Li, T.-J. Mu, and S.-M. Hu, "A large chinese text dataset in the wild," *Journal of Computer Science and Technology*, vol. 34, no. 3, pp. 509–521, 2019. [Online]. Available: https://jcst.ict.ac.cn/en/article/doi/10.1007/s11390-019-1923-y

[92] B. Shi, C. Yao, M. Liao, M. Yang, P. Xu, L. Cui, S. Belongie, S. Lu, and X. Bai, "Icdar2017 competition on reading chinese text in the wild (rctw-17)," in *Proceedings of International Conference on Document Analysis and Recognition*, vol. 1. IEEE, 2017, pp. 1429–1434.

[93] X. Liu, R. Zhang, Y. Zhou, Q. Jiang, Q. Song, N. Li, K. Zhou, L. Wang, D. Wang, M. Liao, M. Yang, X. Bai, B. Shi, D. Karatzas, S. Lu, and C. V. Jawahar, "Icdar 2019 robust reading challenge on reading chinese text on signboard," 2019. [Online]. Available: https://arxiv.org/abs/1912.09641

[94] Y. Sun, J. Liu, W. Liu, J. Han, E. Ding, and J. Liu, "Chinese street view text: Large-scale chinese text reading with partially supervised learning," 2020. [Online]. Available: https://arxiv.org/abs/1909.07808

[95] H. Cheng, P. Zhang, S. Wu, J. Zhang, Q. Zhu, Z. Xie, J. Li, K. Ding, and L. Jin, "M$^6$doc: A large-scale multi-format, multi-type, multi-layout, multi-language, multi-annotation category dataset for modern document layout analysis," 2023. [Online]. Available: https://arxiv.org/abs/2305.08719

[96] B. Deka, Z. Huang, C. Franzen, J. Hibschman, D. Afergan, Y. Li, J. Nichols, and R. Kumar, "Rico: A mobile app dataset for building data-driven design applications," in *Proceedings of the 30th annual ACM symposium on user interface software and technology*, 2017, pp. 845–854.

[97] G. Jaume, H. K. Ekenel, and J.-P. Thiran, "Funsd: A dataset for form understanding in noisy scanned documents," in *Proceedings of International Conference on Document Analysis and Recognition Workshops*, vol. 2.   IEEE, 2019, pp. 1–6.

[98] Z. Huang, K. Chen, J. He, X. Bai, D. Karatzas, S. Lu, and C. Jawahar, "Icdar2019 competition on scanned receipt ocr and information extraction," in *Proceedings of International Conference on Document Analysis and Recognition*.   IEEE, 2019, pp. 1516–1520.

[99] J. Kuang, W. Hua, D. Liang, M. Yang, D. Jiang, B. Ren, and X. Bai, "Visual information extraction in the wild: practical dataset and end-to-end solution," in *Proceedings of International Conference on Document Analysis and Recognition*.   Springer, 2023, pp. 36–53.

[100] Y. Xu, T. Lv, L. Cui, G. Wang, Y. Lu, D. Florencio, C. Zhang, and F. Wei, "XFUND: A benchmark dataset for multilingual visually rich form understanding," in *Proceedings of Annual Meeting of the Association for Computational Linguistics*.   Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 3214–3224. [Online]. Available: https://aclanthology.org/2022.findings-acl.253

[101] W. Yu, C. Zhang, H. Cao, W. Hua, B. Li, H. Chen, M. Liu, M. Chen, J. Kuang, M. Cheng, Y. Du, S. Feng, X. Hu, P. Lyu, K. Yao, Y. Yu, Y. Liu, W. Che, E. Ding, C.-L. Liu, J. Luo, S. Yan, M. Zhang, D. Karatzas, X. Sun, J. Wang, and X. Bai, "Icdar 2023 competition on structured text extraction from visually-rich document images," 2023. [Online]. Available: https://arxiv.org/abs/2306.03287

[102] L. Rujiao, W. Wen, X. Nan, G. Feiyu, Y. Zhibo, W. Yongpan, and X. Gui-Song, "Parsing table structures in the wild," in *Proceedings of IEEE/CVF International Conference on Computer Vision*, October 2021.

[103] F. Yang, L. Hu, X. Liu, S. Huang, and Z. Gu, "A large-scale dataset for end-to-end table recognition in the wild," *Scientific Data*, vol. 10, no. 1, p. 110, 2023.

[104] Y. Liang, Y. Zhang, C. Ma, Z. Zhang, Y. Zhao, L. Xiang, C. Zong, and Y. Zhou, "Document image machine translation with dynamic multi-pre-trained models assembling," in *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2024, pp. 7077–7088.

[105] M. Mathew, D. Karatzas, and C. V. Jawahar, "Docvqa: A dataset for VQA on document images," in *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*. IEEE, 2021, pp. 2199–2208.

[106] Y. Yuan, X. Liu, W. Dikubab, H. Liu, Z. Ji, Z. Wu, and X. Bai, "Syntax-aware network for handwritten mathematical expression recognition," *arXiv preprint arXiv:2203.01601*, 2022.

[107] K. Wang, J. Pan, W. Shi, Z. Lu, M. Zhan, and H. Li, "Measuring multimodal mathematical reasoning with math-vision dataset," *CoRR*, vol. abs/2402.14804, 2024.

[108] W. Yang, Z. Li, D. Peng, L. Jin, M. He, and C. Yao, "Read ten lines at one glance: Line-aware semi-autoregressive transformer for multi-line handwritten mathematical expression recognition," in *Proceedings of the ACM International Conference on Multimedia*, A. El-Saddik, T. Mei, R. Cucchiara, M. Bertini, D. P. T. Vallejo, P. K. Atrey, and M. S. Hossain, Eds. ACM, 2023, pp. 2066–2077.

[109] P. Gervais, A. Fadeeva, and A. Maksai, "Mathwriting: A dataset for handwritten mathematical expression recognition," *CoRR*, vol. abs/2404.10690, 2024.

[110] L. Ding, M. Zhao, F. Yin, S. Zeng, and C.-L. Liu, "A large-scale database for chemical structure recognition and preliminary evaluation," in *Proceedings of the International Conference on Pattern Recognition*, 2022, pp. 1464–1470.

[111] D. Saxton, E. Grefenstette, F. Hill, and P. Kohli, "Analysing mathematical reasoning abilities of neural models," in *Proceedings of the International Conference on Learning Representations*. OpenReview.net, 2019.

[112] R. Zhang, D. Jiang, Y. Zhang, H. Lin, Z. Guo, P. Qiu, A. Zhou, P. Lu, K. Chang, Y. Qiao, P. Gao, and H. Li, "MATHVERSE: does your multi-modal LLM truly see the diagrams in visual math problems?" in *Proceedings of European Conference on Computer Vision*, ser. Lecture Notes in Computer Science, A. Leonardis, E. Ricci, S. Roth, O. Russakovsky, T. Sattler, and G. Varol, Eds., vol. 15066. Springer, 2024, pp. 169–186.

[113] P. Lu, H. Bansal, T. Xia, J. Liu, C. Li, H. Hajishirzi, H. Cheng, K. Chang, M. Galley, and J. Gao, "Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts," in *Proceedings of the International Conference on Learning Representations*. OpenReview.net, 2024.

[114] A. Mishra, S. Shekhar, A. K. Singh, and A. Chakraborty, "Ocr-vqa: Visual question answering by reading text in images," in *Proceedings of International Conference on Document Analysis and Recognition*. IEEE, 2019, pp. 947–952.

[115] H. Feng, W. Zhou, J. Deng, Y. Wang, and H. Li, "Geometric representation learning for document image rectification," in *Proceedings of European Conference on Computer Vision*. Springer, 2022, pp. 475–492.

[116] K. Kafle, S. Cohen, B. Price, and C. Kanan, "Dvqa: Understanding data visualizations via question answering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018.

[117] N. Methani, P. Ganguly, M. M. Khapra, and P. Kumar, "Plotqa: Reasoning over scientific plots," in *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*, March 2020.

[118] M. Mathew, V. Bagal, R. P. Tito, D. Karatzas, E. Valveny, and C. V. Jawahar, "Infographicvqa," *CoRR*, vol. abs/2104.12756, 2021. [Online]. Available: https://arxiv.org/abs/2104.12756

[119] X. Zhong, E. ShafieiBavani, and A. J. Yepes, "Image-based table recognition: data, model, and evaluation," *arXiv preprint arXiv:1911.10683*, 2019.

[120] P. Pasupat and P. Liang, "Compositional semantic parsing on semi-structured tables," in *Proceedings of Annual Meeting of the Association for Computational Linguistics*, C. Zong and M. Strube, Eds. Beijing, China: Association for Computational Linguistics, Jul. 2015, pp. 1470–1480. [Online]. Available: https://aclanthology.org/P15-1142

[121] S. Park, S. Shin, B. Lee, J. Lee, J. Surh, M. Seo, and H. Lee, "Cord: a consolidated receipt dataset for post-ocr parsing," in *Advances in Neural Information Processing Systems Workshop*, 2019.

[122] X. Chen, Z. Zhao, L. Chen, D. Zhang, J. Ji, A. Luo, Y. Xiong, and K. Yu, "Websrc: A dataset for web-based structural reading comprehension," *arXiv preprint arXiv:2101.09465*, 2021.

[123] X. Zhong, J. Tang, and A. Jimeno-Yepes, "Publaynet: Largest dataset ever for document layout analysis," in *Proceedings of International Conference on Document Analysis and Recognition*. IEEE, 2019, pp. 1015–1022.

[124] A. W. Harley, A. Ufkes, and K. G. Derpanis, "Evaluation of deep convolutional nets for document image classification and retrieval," in *Proceedings of International Conference on Document Analysis and Recognition*, 2015.

[125] G. Baechler, S. Sunkara, M. Wang, F. Zubach, H. Mansoor, V. Etter, V. Cărbune, J. Lin, J. Chen, and A. Sharma, "Screenai: A vision-language model for ui and infographics understanding," 2024.

[126] R. Tanaka, K. Nishida, K. Nishida, T. Hasegawa, I. Saito, and K. Saito, "Slidevqa: A dataset for document visual question answering on multiple images," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2023.

[127] A. Kembhavi, M. Salvato, E. Kolve, M. Seo, H. Hajishirzi, and A. Farhadi, "A diagram is worth a dozen images," in *Proceedings of European Conference on Computer Vision*. Springer, 2016, pp. 235–251.

[128] A. Kembhavi, M. Seo, D. Schwenk, J. Choi, A. Farhadi, and H. Hajishirzi, "Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4999–5007.

[129] S. Lee, B. Lai, F. Ryan, B. Boote, and J. M. Rehg, "Modeling multimodal social interactions: New challenges and baselines with densely aligned representations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 14 585–14 595.

[130] G. Zhang, X. Du, B. Chen, Y. Liang, T. Luo, T. Zheng, K. Zhu, Y. Cheng, C. Xu, S. Guo, H. Zhang, X. Qu, J. Wang, R. Yuan, Y. Li, Z. Wang, Y. Liu, Y.-H. Tsai, F. Zhang, C. Lin, W. Huang, and J. Fu, "Cmmmu: A chinese massive multi-discipline multimodal understanding benchmark," 2024. [Online]. Available: https://arxiv.org/abs/2401.11944

[131] P. Lu, S. Mishra, T. Xia, L. Qiu, K.-W. Chang, S.-C. Zhu, O. Tafjord, P. Clark, and A. Kalyan, "Learn to explain: Multimodal reasoning via thought chains for science question answering," *Advances in Neural Information Processing Systems*, vol. 35, pp. 2507–2521, 2022.

[132] X. Yue, T. Zheng, Y. Ni, Y. Wang, K. Zhang, S. Tong, Y. Sun, B. Yu, G. Zhang, H. Sun *et al.*, "Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark," *arXiv preprint arXiv:2409.02813*, 2024.

[133] Q. Jia, X. Yue, S. Huang, Z. Qin, Y. Liu, B. Y. Lin, and Y. You, "Visual perception in text strings," *arXiv preprint arXiv:2410.01733*, 2024.

[134] C. Yao, X. Bai, W. Liu, Y. Ma, and Z. Tu, "Detecting texts of arbitrary orientations in natural images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 1083–1090.

[135] M. He, Y. Liu, Z. Yang, S. Zhang, C. Luo, F. Gao, Q. Zheng, Y. Wang, X. Zhang, and L. Jin, "Icpr2018 contest on robust reading for multi-type web images," in *Proceedings of the International Conference on Pattern Recognition*, 2018, pp. 7–12.

[136] B. Shi, C. Yao, M. Liao, M. Yang, P. Xu, L. Cui, S. Belongie, S. Lu, and X. Bai, "Icdar2017 competition on reading chinese text in the wild (rctw-17)," 2018. [Online]. Available: https://arxiv.org/abs/1708.09585

[137] J. Tang, Z. Yang, Y. Wang, Q. Zheng, Y. Xu, and X. Bai, "Seglink++: Detecting dense and arbitrary-shaped scene text by instance-aware component grouping," *Pattern Recognition*, vol. 96, p. 106954, 2019. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0031320319302511

[138] C. K. Chng, Y. Liu, Y. Sun, C. C. Ng, C. Luo, Z. Ni, C. Fang, S. Zhang, J. Han, E. Ding *et al.*, "Icdar2019 robust reading challenge on arbitrary-shaped text-rrc-art," in *Proceedings of International Conference on Document Analysis and Recognition*. IEEE, 2019, pp. 1571–1576.

[139] X. Zheng, D. Burdick, L. Popa, X. Zhong, and N. X. R. Wang, "Global table extractor (GTE): A framework for joint table identification and cell structure recognition using visual context," in *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*. IEEE, 2021, pp. 697–706. [Online]. Available: https://doi.org/10.1109/WACV48630.2021.00074

[140] X. Dong, P. Zhang, Y. Zang, Y. Cao, B. Wang, L. Ouyang, S. Zhang, H. Duan, W. Zhang, Y. Li *et al.*, "Internlm-xcomposer2-4khd: A pioneering large vision-language model handling resolutions from 336 pixels to 4k hd," *arXiv preprint arXiv:2404.06512*, 2024.

[141] Q. Sun, Y. Cui, X. Zhang, F. Zhang, Q. Yu, Y. Wang, Y. Rao, J. Liu, T. Huang, and X. Wang, "Generative multimodal models are in-context learners," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 14 398–14 409.

[142] J. Ye, H. Xu, H. Liu, A. Hu, M. Yan, Q. Qian, J. Zhang, F. Huang, and J. Zhou, "mplug-owl3: Towards long image-sequence understanding in multi-modal large language models," *arXiv preprint arXiv:2408.04840*, 2024.

[143] W. Wang, Q. Lv, W. Yu, W. Hong, J. Qi, Y. Wang, J. Ji, Z. Yang, L. Zhao, X. Song *et al.*, "Cogvlm: Visual expert for pretrained language models," *arXiv preprint arXiv:2311.03079*, 2023.

[144] Z. Chen, J. Wu, W. Wang, W. Su, G. Chen, S. Xing, M. Zhong, Q. Zhang, X. Zhu, L. Lu *et al.*, "Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 24 185–24 198.

[145] H. Lu, W. Liu, B. Zhang, B. Wang, K. Dong, B. Liu, J. Sun, T. Ren, Z. Li, H. Yang *et al.*, "Deepseek-vl: towards real-world vision-language understanding," *arXiv preprint arXiv:2403.05525*, 2024.

[146] Z. Liu, L. Zhu, B. Shi, Z. Zhang, Y. Lou, S. Yang, H. Xi, S. Cao, Y. Gu, D. Li *et al.*, "Nvila: Efficient frontier visual language models," *arXiv preprint arXiv:2412.04468*, 2024.

[147] A. Hu, H. Xu, L. Zhang, J. Ye, M. Yan, J. Zhang, Q. Jin, F. Huang, and J. Zhou, "mplug-docowl2: High-resolution compressing for ocr-free multi-page document understanding," *arXiv preprint arXiv:2409.03420*, 2024.

[148] A. Young, B. Chen, C. Li, C. Huang, G. Zhang, G. Zhang, H. Li, J. Zhu, J. Chen, J. Chang *et al.*, "Yi: Open foundation models by 01. ai," *arXiv preprint arXiv:2403.04652*, 2024.

[149] C. Wu, X. Chen, Z. Wu, Y. Ma, X. Liu, Z. Pan, W. Liu, Z. Xie, X. Yu, C. Ruan *et al.*, "Janus: Decoupling visual encoding for unified multimodal understanding and generation," *arXiv preprint arXiv:2410.13848*, 2024.

[150] M. Shi, F. Liu, S. Wang, S. Liao, S. Radhakrishnan, D.-A. Huang, H. Yin, K. Sapra, Y. Yacoob, H. Shi *et al.*, "Eagle: Exploring the design space for multimodal llms with mixture of encoders," *arXiv preprint arXiv:2408.15998*, 2024.

[151] H. Laurençon, A. Marafioti, V. Sanh, and L. Tronchon, "Building and better understanding vision-language models: insights and future directions," in *Workshop on Responsibly Building the Next Generation of Multimodal Foundational Models*, 2024.

[152] A. Abouelenin, A. Ashfaq, A. Atkinson, H. Awadalla, N. Bach, J. Bao, A. Benhaim, M. Cai, V. Chaudhary, C. Chen *et al.*, "Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras," *arXiv preprint arXiv:2503.01743*, 2025.

[153] H. Duan, J. Yang, Y. Qiao, X. Fang, L. Chen, Y. Liu, X. Dong, Y. Zang, P. Zhang, J. Wang *et al.*, "Vlmevalkit: An open-source toolkit for evaluating large multi-modality models," in *Proceedings of the ACM International Conference on Multimedia*, 2024, pp. 11 198–11 201.

[154] K. Team, A. Du, B. Yin, B. Xing, B. Qu, B. Wang, C. Chen, C. Zhang, C. Du, C. Wei, C. Wang, D. Zhang, D. Du, D. Wang, E. Yuan, E. Lu, F. Li, F. Sung, G. Wei, G. Lai, H. Zhu, H. Ding, H. Hu, H. Yang, H. Zhang, H. Wu, H. Yao, H. Lu, H. Wang, H. Gao, H. Zheng, J. Li, J. Su, J. Wang, J. Deng, J. Qiu, J. Xie, J. Wang, J. Liu, J. Yan, K. Ouyang, L. Chen, L. Sui, L. Yu, M. Dong, M. Dong, N. Xu, P. Cheng, Q. Gu, R. Zhou, S. Liu, S. Cao, T. Yu, T. Song, T. Bai, W. Song, W. He, W. Huang, W. Xu, X. Yuan, X. Yao, X. Wu, X. Zu, X. Zhou, X. Wang, Y. Charles, Y. Zhong, Y. Li, Y. Hu, Y. Chen, Y. Wang, Y. Liu, Y. Miao, Y. Qin, Y. Chen, Y. Bao, Y. Wang, Y. Kang, Y. Liu, Y. Du, Y. Wu, Y. Wang, Y. Yan, Z. Zhou, Z. Li, Z. Jiang, Z. Zhang, Z. Yang, Z. Huang, Z. Huang, Z. Zhao, and Z. Chen, "Kimi-VL technical report," 2025. [Online]. Available: https://arxiv.org/abs/2504.07491

# NeurIPS Paper Checklist

1. **Claims**

   Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

   Answer: [Yes]

   Justification: We have clarified the motivation and contributions of our work.

   Guidelines:

   - The answer NA means that the abstract and introduction do not include the claims made in the paper.
   - The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
   - The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
   - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. **Limitations**

   Question: Does the paper discuss the limitations of the work performed by the authors?

   Answer: [Yes]

   Justification: Due to space limitations, we provide the limitation section in the Appendix.

   Guidelines:

   - The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
   - The authors are encouraged to create a separate "Limitations" section in their paper.
   - The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
   - The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
   - The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
   - The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
   - If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
   - While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. **Theory assumptions and proofs**

   Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

   Answer: [NA]

Justification: The evaluation metric formulas are provided in the Appendix for clarity. As this work focuses on benchmarking, no additional assumptions or formal proofs are presented.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. **Experimental result reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We have provided the data and the evaluation scripts. The necessary guidance information and the real test example are also provided.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general. releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. **Open access to data and code**

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We have provided the data and the evaluation scripts.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. **Experimental setting/details**

   Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

   Answer: [Yes]

   Justification: We provide the detailed data construction process and the evaluation settings in the paper and the Appendix.

   Guidelines:

   - The answer NA means that the paper does not include experiments.
   - The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
   - The full details can be provided either with the code, in appendix, or as supplemental material.

7. **Experiment statistical significance**

   Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

   Answer: [No]

   Justification: The results are obtained by running inference with publicly available pre-trained models using default parameters.

   Guidelines:

   - The answer NA means that the paper does not include experiments.
   - The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
   - The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. **Experiments compute resources**

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The evaluation experiments were conducted on GPUs, and the detailed evaluation setup is provided in the Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. **Code of ethics**

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We've checked the Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The societal impacts of our benchmark is discussed in the Appendix.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. **Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: All samples included in our benchmark have been manually filtered to ensure that the content is safe for release.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. **Licenses for existing assets**

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Some of the data used in our study are derived from existing academic datasets, with detailed information provided in the Appendix. In addition, we use publicly pre-trained models to evaluate their performance on our benchmark.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New assets**

    Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

    Answer: [Yes]

    Justification: We establish a unified instruction format for 23 text-oriented tasks and provide corresponding evaluation metrics. Additionally, self-annotated data are included in both the public and private test sets of our benchmark.

    Guidelines:
    - The answer NA means that the paper does not release new assets.
    - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
    - The paper should discuss whether and how consent was obtained from people whose asset is used.
    - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

    Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

    Answer: [Yes]

    Justification: Our benchmark construction involved human filtering and annotation. All major contributors to this process are listed as co-authors.

    Guidelines:
    - The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
    - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
    - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

    Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

    Answer: [No]

    Justification: The evaluation process doesn't involve human subjects.

    Guidelines:
    - The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. **Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: LLMs were used solely for language editing and expression refinement. They did not contribute to the core methodology, scientific rigor, or originality of the research.

Guidelines:
- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (`https://neurips.cc/Conferences/2025/LLM`) for what should or should not be described.

# A Technical Appendices and Supplementary Material

This supplementary material contains the following content:

- **Sec. A.1**: Comparison experiments between LMMs and some text-centric expert models.
- **Sec. A.2**: Data collection.
- **Sec. A.3**: Task definitions.
- **Sec. A.4**: Additional statistics of *OCRBench v2*.
- **Sec. A.5**: Evaluation metrics.
- **Sec. A.6**: Experimental setting for the evaluation process.
- **Sec. A.7**: Compute resources for the evaluation process.
- **Sec. A.8**: Evaluation results for LMMs on *OCRBench v2*.
- **Sec. A.9**: Potential factors affecting OCR capabilities
- **Sec. A.10**: Visualization samples for task examples.
- **Sec. A.11**: Visualization samples for failure cases.
- **Sec. A.12**: Biases during the data construction.
- **Sec. A.13**: Discussion of broader impacts.
- **Sec. A.14**: Discussion of limitations.

## A.1 Comparison with LMMs and Text-centric Expert Models

**Comparison with text recognizers.** We compare LMMs with several representative scene text recognizers, including CRNN [67], ABINet [68], ASTER [69], MASTER [70], and SVTR [71], on the text recognition task. The weights of these models are loaded from mmocr[2]. The results are shown in Tab. 8, where we selected 5 representative LMMs, including Qwen2.5VL-7B [66], InternVL3-14B [23], GPT4o [1], Gemini1.5-Pro [60], and Step-1V [62]. The results demonstrate that LMMs exhibit remarkable text recognition capabilities, validating our motivation to evaluate LMMs on more challenging OCR-related tasks.

Table 8: Comparison between LMMs and text recognizers.

| Method | Accuracy |
|---|---|
| CRNN [67] | 38.1 |
| ABINet [68] | 62.4 |
| ASTER [69] | 50.0 |
| MASTER [70] | 54.1 |
| SVTR [71] | 57.8 |
| Qwen2.5VL-7B [66] | 73.0 |
| InternVL3-14B [23] | 71.1 |
| GPT4o [1] | 74.1 |
| Gemini1.5-Pro [60] | 64.1 |
| Step-1V [62] | 75.4 |

**Comparison with text spotters.** We also compare LMMs with ABCNet series [72, 73] and TESTR [74] on the text spotting task. The ABCNet series utilize the official weights[3], and TESTR is also initialized with its publicly released checkpoint[4]. These models were fine-tuned with Total-Text [75]. The results are shown in Tab. 9. Although LMMs demonstrate promising capabilities in text recognition, there remains notable potential for improvement in the text spotting task.

**Comparison with GOT.** We notice a recent work, GOT [76], that can parse the textual elements within images. We conduct comparison experiments between GOT and some representative LMMs, and the results are shown in Tab. 10. We observe that LMMs show advantages in general text recognition, while GOT demonstrates better performance in the document parsing task.

---

[2] https://github.com/open-mmlab/mmocr
[3] https://github.com/aim-uofa/AdelaiDet
[4] https://github.com/mlpc-ucsd/TESTR

Table 9: Comparison between LMMs and text spotters.

| Method | F1 score |
|---|---|
| ABCNet [72] | 32.2 |
| ABCNetV2 [73] | 44.2 |
| TESTR [74] | 51.8 |
| Qwen2.5VL-7B [66] | 1.2 |
| InternVL3-14B [23] | 11.2 |
| Gemini1.5-Pro [60] | 13.5 |
| GPT4o [1] | 0 |
| Step-1V [62] | 7.2 |

Table 10: Comparison between LMMs and GOT [76].

| Method | Rec | FG-Rec | Full-Rec | Doc-Parse |
|---|---|---|---|---|
| GOT [76] | 64.1 | 52.9 | 73.3 | 53.9 |
| Qwen2.5VL-7B [66] | 73.0 | 36.4 | 84.2 | 39.1 |
| InternVL3-14B [23] | 71.1 | 36.4 | 83.0 | 36.9 |
| GPT4o [1] | 74.1 | 13.8 | 54.1 | 35.9 |
| Gemini1.5-Pro [60] | 64.1 | 22.9 | 83.9 | 40.5 |
| Step-1V [62] | 76.8 | 24.8 | 74.8 | 36.0 |

## A.2 Data Collection

**Text Recognition.** The data for text recognition task are sampled from ICDAR2013 [77], SVT [78], IIIT5K [79], ICDAR2015 [80], SCUT-CTW1500 [81], COCO-Text [82], CUTE80 [83], TotalText, SVTP [84], WordArt [85], NonSemanticText [14], IAM [86], ORAND-CAR-2014 [87], HOST [88], and WOST [88]. Meanwhile, CAPTCHA (Completely Automated Public Turing Test to Tell Humans Apart) images are sourced from a CAPTCHA dataset[5] and a number CAPTCHA dataset[6]. Additionally, dot matrix images in the text recognition task are manually collected from the web page.

**Fine-grained Text Recognition.** In the fine-grained text recognition task, images are sampled from the test sets of Fox [17], Totaltext, COCO-Text, CTW1500 [89], and ICDAR2015. We use the original annotations for Fox, while the other datasets are manually re-annotated.

**Full-page OCR.** The data sources for full-page OCR task include Fox, HierText [90], CTW [91], RCTW-17 [92], ReCTS [93], LSVT2019 [94], M6Doc [95], and CDLA[7].

**Text Grounding.** The images for the text grounding task are sampled from testset of Totaltext, COCO-Text, CTW1500, and ICDAR2015. QA pairs and bounding boxes annotations are based on their official OCR annotations.

**VQA with Position.** The images used for VQA with position task are sampled from the test sets of TextVQA [10] and RICO [96], with QA pairs and bounding box annotations derived from their original datasets.

**Text Spotting.** The data sources for the text spotting task include Totaltext, COCO-Text, CTW1500, and ICDAR2015.

**Key Information Extraction.** The data sources for key information extraction task include FUNSD [97], SROIE [98], POIE [99], M6Doc, XFUND [100], ICDAR2023-SVRD [101], and a private dataset of photographed receipts.

**Key Information Mapping.** The data sources for the key information mapping task include FUNSD and POIE.

---

[5] https://aistudio.baidu.com/datasetdetail/159309
[6] https://www.heywhale.com/mw/dataset/5e5e56b6b8dfce002d7ee42c/file
[7] https://github.com/buptlihang/CDLA

**Handwritten Content Extraction.** This task's data is our private data, which contains real exam paper data with student information removed and manually annotated QA pairs.

**Table Parsing.** The images for table parsing task are selected from MMTab [28], WTW [102], TabRecSet [103] and flush table recognition competition[8].

**Chart Parsing.** The data sources for the chart parsing task come from OneChart [44] and MMC [29].

**Document Parsing.** The data sources for document parsing task come from DoTA [104], DocVQA [105], M6Doc, and CDLA.

**Formula Recognition.** The data sources for the formula Recognition task includes HME100K [106], IM2LATEX-100K [107], M2E [108], MathWriting [109], MLHME-38K[9], CASIA-CSDB [110], and some private data.

**Math QA.** The data sources for the math QA task includes MathMatics [111], MathVerse [112], MathVision [107], and MathVista [113].

**Text Counting.** The data for the text counting task are collected from IIIT5K, SVT, ICDAR2013, HierText, and TotalText.

**Cognition VQA.** The data sources for the cognition VQA task include EST-VQA [12], OCRVQA [114], ST-VQA [11], TEXTVQA, DIR300 [115], ChartQA [41], DVQA [116], PlotQA [117], InfoVQA [118], WTW, PubTabNet [119], WTQ [120], CORD [121], LLaVAR [30], WebSRC [122], DocVQA, M6Doc, XFUND, Publaynet [123], RVL-CDIP [124], ScreenQA [125], SlideVQA [126], a movie poster collection dataset[10], a website screenshot collection dataset[11], and a private receipt photograph dataset.

**Diagram QA.** The data sources for the diagram QA task include AI2D [127] and TextBookQA [128].

**Document Classification.** The images for the document classification task are collected from RVL-CDIP.

**Reasoning VQA.** The reasoning VQA task shares some common data sources with the cognition VQA task. Additionally, portions of the reasoning VQA dataset are drawn from MMSI [129] and CMMMU [130].

**Science QA.** The images and annotations of the science QA task are collected from ScienceQA [131] and MMMU-Pro [132]

**APP Agent.** The data source of the APP agent task is RICO.

**ASCII Art Classification.** The data sources for the ASCII art classification task is ASCIIEval [133].

**Text Translation.** The datasets collected for text translation task includes memes[12], MSRA-TD500 [134], MTWI2018 [135], M6Doc, ICDAR2023-SVRD, EST-VQA, RCTW17 [136], DAST1500 [137], XFUND, ArT2019 [138], ChartQA, CDLA, ICDAR2015, SlideVQA, Fintabnet [139], ScienceQA, InfoVQA, COMICS-Dialogue[13], and ExpressExpense SRD[14].


## A.3 Task Definitions

In this section, we introduce the definition of each task, and the visualizations for each task can be found in Sec. A.10.

**Text Recognition.** Text recognition refers to the fundamental OCR ability on text image patches, which asks LMMs to read the text content. To comprehensively evaluate LMMs' text recognition ability across diverse scenarios, our collection incorporates various text types, including regular text,

---

Table 11: The number of images included in each scene category in public data.

| Scene | Number | Scene | Number | Scene | Number |
|---|---|---|---|---|---|
| Schematic diagram | 1238 | Scientific paper | 799 | Word | 728 |
| Table(filled) | 705 | Chart | 620 | Receipts | 609 |
| Questions | 581 | Mathematical formula | 475 | Product labels | 434 |
| Phone screenshot | 431 | Indoor scenes | 395 | Industry research reports | 343 |
| Poster | 264 | Street scene | 224 | ASCII Art | 199 |
| Shop sign | 189 | Financial reports | 153 | Chemical formula | 149 |
| Textbook | 148 | Magazine | 146 | Email | 111 |
| Web screenshot | 99 | Details page | 95 | Verification code | 87 |
| Resumes | 67 | Illustration | 61 | Newspaper | 52 |
| Road signs | 43 | Menus | 31 | Notify | 30 |
| Questionnaire | 29 | | | | |



(a) Overview of the eight tasks in OCRBench v2     (b) The ratios of each type of tasks in OCRBench v2

Figure 5: **Overview of the eight testable text-reading capabilities and associated tasks in OCRBench v2**. Each color represents a distinct capability type.

irregular text, artistic text, handwriting text, digit string text, non-semantic text, occluded text, doc matrix text, and CAPTCHA text.

**Fine-grained Text Recognition.** This task requires LLMs to read and comprehend textual content within the given region. It evaluates LLMs' fine-grained perception capabilities in understanding text in natural scenes and documents.

**Full-page OCR.** Full-page OCR [17] task requires LMMs to extract and recognize all text content from the given images. Converting text into digital format facilitates subsequent processing and analysis of text images.

**Text grounding.** In this task, users would provide a text string and require LMMs to locate its specific location, evaluating LMMs' fine-grained perception capabilities.

**VQA with Position.** For VQA with position task, LMMs need to not only respond to the question but also provide the exact position coordinates that directly correspond to the answer. We ask LMMs to output both information in JSON format for convenient evaluation, and the coordinates are required to be normalized with image sizes and scaled to the range of $[0, 1000]$.

**Text Spotting.** Text spotting task needs LMMs to output the localization and content of all appeared text simultaneously. Due to the interference of background elements and the large number of text instances, this task demands high fine-grained perception capabilities from the model. Besides, the coordinates are required to be normalized with image sizes and scaled to the range of $[0, 1000]$.

Figure 6: The quantity distribution of English tasks of public data.



Figure 7: The quantity distribution of Chinese tasks of public data.

**Key Information Extraction.** The key information extraction task is to extract the necessary information from densely arranged text. In this task, we provide some desired entities as keys and demand LMMs to output the corresponding values to form the output JSON string.

**Key Information Mapping.** In this task, we provide a set of entity keys and their corresponding values in the prompt. The LMMs are then asked to match and pair these keys with their respective values into groups.

**Handwritten Content Extraction.** To investigate the information extraction capabilities of LMMs in educational scenarios, we collect some Chinese examination papers, containing both printed question text and handwritten student responses. There are four types of questions in these examination papers, including single-choice, multiple-choice, true or false, and brief response questions. The prompts require LMMs to extract the handwritten content for specific questions.

**Table Parsing.** Table parsing task requires LMMs to parse the given table into structured text, including Markdown and HTML format.

**Chart Parsing.** Apart from tables, charts can also be converted to structured information. In this task, LLMs are required to transform visual charts into JSON format.

**Document Parsing.** In the document parsing task, both text and the complex elements, including charts, tables, and formulas, are required to be parsed.

**Formula Recognition.** This task asks LMMs to recognize the given formula in the LaTeX format. The collection includes mathematical and chemical formulas.

Figure 8: The OCR lines distribution of English tasks of public data.



Figure 9: The OCR lines distribution of Chinese tasks of public data.

**Math QA.** Math QA task evaluates the LMMs' mathematical calculation ability. In particular, we render the mathematical problem description and related figures into images and ask LMMs to answer the questions within the images.

**Text Counting.** Text counting task is built to evaluate the quantity property perceiving ability of LMMs, including the character frequency in words and the word counting in the given image.

**Cognition VQA.** In *OCRBench v2*, we split text-centric VQA instructions into cognition VQA and Reasoning VQA based on whether the answers can be directly found in the images. Cognition VQA task refers to the instructions where answers are explicitly present in the given image. This task evaluates the fundamental text-centric question-answering ability based on visual content.

**Diagram QA.** In the diagram QA task, LMMs need to respond to the question about the given diagrams, reflecting LMMs' ability to understand the relationship between the visual elements.

**Document Classification.** Document classification task asks LMMs to classify the category of the given document image. The included categories are letters, forms, emails, handwritten documents, advertisements, scientific reports, scientific publications, specifications, file folders, news articles, budgets, invoices, presentations, questionnaires, resumes, and memos.

**Reasoning VQA.** In reasoning VQA tasks, the answers often do not directly appear in the image. This forces LMMs to perform logical reasoning to respond to questions based on visual information.

**Science QA.** In the Science QA task, LMMs are required to respond to the scientific problem. We use PaddleOCR[15] to extract text from the collected images and filter out those with fewer than four OCR

---

[15] https://github.com/PaddlePaddle/PaddleOCR

results. Additionally, when extra subject-related knowledge is provided by the source, we incorporate it by rendering it into the images.

**APP Agent.** For the APP agent task, LMMs need to understand the relationship between textual content, icons, and world knowledge to respond to the question from the user, simulating the real-world application scene.

**ASCII Art Classification.** We incorporate a recent image classification task that uses images composed purely of ASCII characters [133]. This task is included in *OCRBench v2* to evaluate LMMs' ability to assess LMMs' pattern recognition and visual abstraction abilities.

**Text Translation.** In the text translation task, LMMs need to execute translation between Chinese and English texts, evaluating LMMs' semantic understanding abilities.

### A.4 Additional Statistics of OCRBench v2

**Scene Coverage.** Our dataset can be divided into 31 classic scenes according to the scene of the image. The specific scenes and the corresponding number of pictures are shown in Tab. 11.

**Statistics of each task.** Fig. 5 shows an overview of each task in *OCRBench v2*.The distribution of 23 tasks in *OCRBench v2* is displayed in Fig. 6 and Fig. 7. Additionally, we calculate and present the average number of OCR text lines per task in Fig. 8 and Fig. 9. As illustrated in these figures, the task distribution is well-balanced, with each task containing adequate textual information for analysis.

### A.5 Evaluation Metrics

**Parsing Type.** We use Tree-Edit-Distance-based Similarity (TEDS) [63] to evaluate parsing tasks, which require LMMs to transform the images to structured formats. Tree Edit Distance (TED) refers to the minimum number of edits to transform one tree into another. TEDS is based on TED to calculate the similarity of two trees. Assuming $T_1$ and $T_2$ are two different trees, $TED(T_1, T_2)$ refers to their TED, and the TEDS is defined as:

$$TEDS(T_1, T_2) = 1 - \frac{TED(T_1, T_2)}{\max(|T_1|, |T_2|)}, \tag{1}$$

where $|T_1|$ and $|T_2|$ is the number of nodes of trees, $TED(T_1, T_2)$ can be calculated by dynamic programming algorithm. If $T_1$ and $T_2$ are identical, then their TEDS equals 1. As the structural difference between two trees increases, their TED value becomes larger, resulting in the TEDS approaching 0.

**Localization Type.** In the text referring and spotting tasks, LMMs are required to provide regression bounding boxes of target objects. IoU score is adopted to measure the distance between the predicted regions and the ground truth.

$$IoU(B_1, B_2) = \frac{Intersect(B_1, B_2)}{Union(B_1, B_2)}, \tag{2}$$

where $Intersect(B_1, B_2)$ refers to the overlap area of bounding box $B_1$ and $B_2$, while $Union(B_1, B_2)$ refers to their union area.

**Extraction Type.** The F1 score is used to evaluate LMMs' relation extraction capability. Given the predicted and ground truth Key-Value pairs, the F1 score is formulated as follows:

$$Precision = \frac{N_3}{N_2}, \tag{3}$$

$$Recall = \frac{N_3}{N_1}, \tag{4}$$

$$Fmean = \frac{2 * Precision * Recall}{Precision + Recall}, \tag{5}$$

where $N_1$, $N_2$, and $N_3$ denote the number of ground-truth Key-Value pairs, predicted Key-Value pairs, and correctly matched Key-Value pairs, respectively.

**Long Reading Type.** To evaluate LMMs' ability to recognize text across entire paragraphs or pages, BLEU [64], METEOR [65], F1 score, and normalized edit distance are employed. And the final score is the average value of these metrics.

BLEU evaluates prediction quality by comparing n-gram match rates between the prediction and ground truth sequences. For each n-gram type, precision is calculated as the ratio of matching n-grams to total predicted n-grams. The final BLEU score is the geometric mean of these precision values multiplied by a penalty $BP$, which is defined as:

$$BLEU = BP * exp(\sum_{n=1}^{N} w_n \log p_n), \tag{6}$$

$$BP = \begin{cases} 1 & L_p \geq L_g \\ e^{(1-\frac{L_p}{L_g})} & L_p < L_g \end{cases}, \tag{7}$$

where $p_n$ represents the precision of n-grams, $L_p$ represents the length of prediction sequence, $L_g$ represents the length of ground truth sequence, $w_n$ is weight factor, usually evenly distributed ($w_n = \frac{1}{N}$). Typically, $N$ is set to 4.

METEOR employs a semantic-aware matching strategy with four levels. 1) Exact Match: words in the prediction that are identical to the ground truth. 2) Stem match: matching words that have the same word stem. 3) Synonym Match: matching words based on synonymous relationships. 4) Paraphrase Match: Matching similar phrases at the phrase level. These matches are combined to calculate precision and recall, from which a weighted harmonic mean F1 score is derived as:

$$P_{meteor} = \frac{N_{match}}{N_{pred}}, \tag{8}$$

$$R_{meteor} = \frac{N_{match}}{N_{gt}}, \tag{9}$$

$$F_{meteor} = \frac{10 * P_{meteor} * R_{meteor}}{P_{meteor} + 9 * R_{meteor}}, \tag{10}$$

where $N_{match}$, $N_{pred}$, and $N_{gt}$ represent the number of matched items, words in prediction, and words in ground truth, respectively. The final METEOR score is obtained by multiplying the $F_{meteor}$ by the penalty adjustment factor. The calculation is formulated as follows:

$$METEOR = F_{meteor} * (1 - BP_{meteor}), \tag{11}$$

$$BP_{meteor} = 0.5 * \frac{N_{chunk}}{N_{match}}, \tag{12}$$

where $N_{chunk}$ refers to the number of contiguous matching phrases. More chunks indicate greater word order differences, resulting in a heavier penalty.

The calculation method of the F1 score in long reading metrics follows the same approach as discussed in extraction metrics, as shown in Equations 3, 4, 5.

Normalized Edit Distance (NED) measures string similarity by computing the minimum number of operations needed to transform one string into another. And then NED is normalized by the length of the longer string. The calculation is formulated as follows:

$$NED(S_1, S_2) = \frac{ED(S_1, S_2)}{\max(len(S_1), len(S_2))} \tag{13}$$

where $ED(S_1, S_2)$ represents the edit distance between the prediction string $S_1$ and the ground truth $S_2$. The $NED$ value of 0 indicates identical strings, while 1 indicates completely different strings.

**Counting Type.** In *OCRBench v2*, character frequency counting and word counting tasks are included. For character frequency, we use exact match evaluation since the answers are typically single-digit

integers. For word counting, we evaluate using the L1 distance between predicted and ground truth counts, normalized to $[0, 1]$ based on the ground truth. This can be formulated as follows:

$$score = \begin{cases} 0 & C_{pred} \leq 0 \\ 1 - \frac{|C_{pred} - C_{gt}|}{C_{gt}}) & 0 < C_{pred} < 2 * C_{gt} \\ 0 & C_{pred} \geq 2 * C_{gt} \end{cases}, \tag{14}$$

where $C_{pred}$ and $C_{gt}$ denote the predicted count and ground truth count, respectively.

**Basic VQA Type.** The remaining tasks in *OCRBench v2* are basic VQA types, and we employ different evaluation metrics based on question types. For multiple-choice questions, we use exact matching between predictions and answer options. In other cases, we check whether the ground truth is contained in the prediction for answers shorter than 5 words, and use ANLS for longer answers.

### A.6 Experimental setting

The detailed public data construction are shown in Sec. A.2 and Sec. A.5. Private data consists of unlabeled images collected manually from websites and real life. At the same time, we annotated and checked the private test set to ensure the quality. The environment configuration of each open-source model experiment strictly complies with the official version and uses the official pre-trained model and inference code. The model parameters of the open-source model and the API parameters of the closed-source model use the official default parameters for fair. Specifically, we use the official API versions: GPT-4o (gpt-4o-2024-08-06), GPT-4o-mini (gpt-4o-mini-2024-07-18), and Gemini 1.5 Pro (gemini1.5-pro-002).

### A.7 Compute resources

Evaluations of open-source models were conducted on 8×NVIDIA GeForce RTX 4090 (24GB) and a NVIDIA H800 Tensor Core GPU (80GB). The closed-source experiments obtained the results by calling the official API.

### A.8 Results and Discussions

Tab. 12, Tab. 13, Tab. 14, and Tab. 15 exhibit the results of 39 open-source models and 5 closed-source models on the public and private test sets of *OCRBench v2*

**Evaluation results on public data** are shown in Tab. 12 and Tab. 13. Most LMMs performed well in tasks such as Understanding, Recognition, Extraction, which shows that current models have basic OCR capabilities. However, they performed poorly in tasks such as Referring, Spotting, Parsing, and Calculation. The scores of all models are basically below 50 points, which shows that the models still lack the ability in text localization, logical reasoning, and understanding complex elements.

**Evaluation results on private data** are shown in Tab. 14 and Tab. 15. The performance trends of the models on private and public datasets are consistent. In addition, most models perform worse on private datasets than on public datasets, which shows that private data may be more challenging for LMMs due to the lack of training, and also reflects the importance of private data construction.

### A.9 Potential Factors Affecting OCR Capabilities

**High-Res Visual Encoders.** Since text often appears small in images, the resolution setting of the visual encoder could be a key factor affecting the text perception ability [51]. Here we change the input resolution of the LMMs and observe the performance changes. In particular, InternVL2-8B is chosen, and the resolution setting includes $448$, $896$, and dynamic. Tab. 16 lists the results. Indeed, when the input resolution increases from $448$ to $896$, the performance increases by $4.1\%$.

**Pre-provided OCR Information.** To study the impact of OCR information, we use PaddleOCR[16] to pre-extract OCR results and incorporate them with prompts. Tab. 17 shows the results. We observe

---

[16]`https://github.com/PaddlePaddle/PaddleOCR`

Table 12: **Evaluation of existing LMMs on English tasks of OCRBench v2's public data**. "Recognition", "Referring", "Spotting", "Extraction", "Parsing", "Calculation", "Understanding", and "Reasoning" refer to text recognition, text referring, text spotting, relation extraction, element parsing, mathematical calculation, visual text understanding, and knowledge reasoning, respectively. Higher values indicate better performance. Best performance is in boldface, and the second best is underlined. The notations apply to all subsequent figures.

| Method | Recognition | Referring | Spotting | Extraction | Parsing | Calculation | Understanding | Reasoning | Average |
|---|---|---|---|---|---|---|---|---|---|
| Open-source LMMs | | | | | | | | | |
| LLaVA-Next-8B [49] | 41.3 | 18.8 | 0 | 49.5 | 21.2 | 17.3 | 55.2 | 48.9 | 31.5 |
| LLaVA-OV-7B [50] | 46.0 | 20.8 | 0.1 | 58.3 | 25.3 | 23.3 | 64.4 | 53.0 | 36.4 |
| Monkey [51] | 35.2 | 0 | 0 | 16.6 | 16.3 | 14.4 | 59.8 | 42.3 | 23.1 |
| TextMonkey [7] | 39.1 | 0.7 | 0 | 19.0 | 12.2 | 19.0 | 61.1 | 40.2 | 23.9 |
| XComposer2-4KHD [140] | 45.1 | 21.8 | 0.1 | 15.9 | 11.7 | 15.7 | 66.8 | 45.9 | 27.9 |
| Molmo-7B [52] | 52.4 | 21.3 | 0.1 | 45.5 | 7.6 | 28.5 | 65.3 | 55.0 | 34.5 |
| Cambrian-1-8B [53] | 45.3 | 21.5 | 0 | 53.6 | 19.2 | 19.5 | 63.5 | 55.5 | 34.7 |
| Pixtral-12B [54] | 48.9 | 21.6 | 0 | 66.3 | 35.5 | 29.8 | 66.9 | 53.7 | 40.3 |
| EMU2-chat [141] | 42.1 | 0.2 | 0 | 12.5 | 8.1 | 11.2 | 42.7 | 33.4 | 18.8 |
| mPLUG-Owl3 [142] | 41.6 | 14.0 | 0.6 | 24.4 | 10.9 | 11.1 | 52.2 | 46.0 | 25.1 |
| CogVLM-chat [143] | 50.9 | 0 | 0 | 0.2 | 8.4 | 15.0 | 58.1 | 41.7 | 21.8 |
| Qwen-VL [4] | 34.6 | 7.5 | 0 | 18.2 | 20.0 | 8.1 | 57.2 | 41.1 | 23.3 |
| Qwen-VL-chat [4] | 34.5 | 4.1 | 0 | 25.9 | 14.0 | 13.8 | 55.7 | 39.5 | 23.4 |
| Qwen2-Vl-7B [66] | <u>72.1</u> | **47.9** | **17.5** | 82.5 | 25.5 | 25.4 | 78.4 | 61.5 | 51.4 |
| Qwen2.5-VL-7B [21] | 68.8 | 25.7 | 1.2 | 80.2 | 30.4 | 38.2 | 73.2 | 56.2 | 46.7 |
| InternVL2-8B [144] | 49.9 | 23.1 | 0.5 | 65.2 | 24.8 | 26.7 | 73.5 | 52.9 | 39.6 |
| InternVL2-26B [144] | 63.4 | 26.1 | 0 | 76.8 | 37.8 | 32.3 | **79.4** | 58.9 | 46.8 |
| InternVL2.5-8B [23] | 59.0 | 25.0 | 1.4 | 77.5 | 35.1 | 29.4 | 75.3 | 57.2 | 45.0 |
| InternVL2.5-26B [23] | 65.6 | 26.1 | 1.6 | <u>86.9</u> | 36.2 | 37.4 | 78.3 | **62.9** | 49.4 |
| InternVL3-8B [23] | 68.6 | 30.4 | 8.8 | <u>85.3</u> | 34.0 | 27.1 | 77.5 | 60.3 | 49.0 |
| InternVL3-14B [23] | 67.3 | 36.9 | 11.2 | **89.0** | 38.4 | 38.4 | <u>79.2</u> | 60.5 | **52.6** |
| Deepseek-VL-7B [145] | 37.1 | 15.4 | 0 | 23.5 | 14.6 | 20.8 | 53.3 | 52.9 | 27.2 |
| Deepseek-VL2-Small [55] | 62.7 | 28.0 | 0.1 | 77.5 | 32.7 | 14.3 | 77.1 | 53.9 | 43.3 |
| MiniCPM-V-2.6 [56] | 66.8 | 6.0 | 0.8 | 62.0 | 28.8 | 32.4 | 73.7 | 52.1 | 40.3 |
| MiniCPM-o-2.6 [56] | 66.9 | 29.5 | 0.5 | 70.8 | 33.4 | 31.9 | 69.9 | 57.9 | 45.1 |
| GLM-4V-9B [57] | 61.8 | 22.6 | 0 | 71.7 | 31.6 | 22.6 | 72.1 | 58.4 | 42.6 |
| VILA1.5-8B [146] | 35.3 | 15.5 | 0 | 21.1 | 12.7 | 17.3 | 46.3 | 40.3 | 23.6 |
| LLaVAR [30] | 37.3 | 0 | 0 | 1.0 | 9.9 | 12.3 | 34.6 | 27.0 | 15.3 |
| UReader [33] | 22.4 | 0.1 | 0 | 0 | 9.2 | 7.9 | 41.0 | 29.1 | 13.7 |
| DocOwl2 [147] | 24.0 | 9.7 | 0 | 13.4 | 13.5 | 8.8 | 53.7 | 32.0 | 19.4 |
| Yi-VL-6B [148] | 28.9 | 2.9 | 0 | 9.7 | 12.9 | 15.8 | 36.1 | 32.0 | 17.3 |
| Janus-1.3B [149] | 46.1 | 0 | 0 | 0.2 | 14.5 | 13.5 | 36.0 | 39.1 | 18.7 |
| Eagle-X5-7B [150] | 34.7 | 17.8 | 0 | 21.7 | 20.6 | 21.5 | 61.0 | 42.6 | 27.5 |
| Idefics3-8B [151] | 23.8 | 13.2 | 0 | 63.2 | 23.8 | 23.0 | 65.8 | 44.9 | 32.2 |
| Phi-4-MultiModal [152] | 63.7 | 16.4 | 0 | 40.4 | 19.1 | 18.3 | 69.8 | 53.9 | 35.2 |
| SAIL-VL-1.6-8B [153] | 67.7 | 28.6 | 2.8 | 70.5 | 25.9 | 29.5 | 73.9 | 59.7 | 44.8 |
| Kimi-VL-A3B-16B [154] | 56.5 | 13.8 | 0 | 59.2 | 33.8 | 32.9 | 75.5 | 56.7 | 41.1 |
| Ovis1.6-3B [58] | 59.2 | 14.3 | 0 | 65.0 | 32.1 | 29.0 | 69.8 | 56.8 | 40.8 |
| Ovis2-8B [58] | **73.2** | 24.6 | 0.7 | 62.4 | **44.8** | 40.6 | 72.7 | <u>62.6</u> | 47.7 |
| Closed-source LMMs | | | | | | | | | |
| GPT-4o [1] | 61.2 | 26.7 | 0 | 77.5 | 36.3 | <u>43.4</u> | 71.1 | 55.5 | 46.5 |
| GPT-4o-mini [59] | 57.9 | 23.3 | 0.6 | 70.8 | 31.5 | 38.8 | 65.9 | 55.1 | 43.0 |
| Gemini-Pro [60] | 61.2 | <u>39.5</u> | <u>13.5</u> | 79.3 | <u>39.2</u> | **47.7** | 75.5 | 59.3 | <u>51.9</u> |
| Claude3.5-sonnet [61] | 62.2 | 28.4 | 1.3 | 56.6 | 37.8 | 40.8 | 73.5 | 60.9 | 45.2 |
| Step-1V [62] | 67.8 | 31.3 | 7.2 | 73.6 | 37.2 | 27.8 | 69.8 | 58.6 | 46.7 |

that adding OCR information does not help much. This suggests that *OCRBench v2* evaluates LMMs capabilities across multiple dimensions, rather than solely focusing on text recognition abilities.

**Connection Between OCR and LLMs.** We further explore a direct pipeline by first extracting OCR information and then by feeding it directly into Qwen2.5. Unlike LMMs, this pipeline separates OCR and language modeling into distinct stages. The results shown in Tab. 17 suggest that Qwen2-VL-7B outperforms Qwen2.5 with OCR information, demonstrating LMMs' remarkable ability to incorporate both textual and visual features efficiently.

## A.10    Samples for Each Task

As show in Fig. 10 to Fig. 18 , there are 23 OCR tasks included in *OCRBench v2*. Among them, Fig. 10 to Fig. 16 present examples of English tasks, including text recognition, diagram QA, text counting, formula recognition, math QA, VQA with position, ASCII art classification, reasoning VQA, text translation, APP agent, table parsing, cognition VQA, document classification, science QA, chart parsing, key information extraction, full-page OCR, text spotting, fine-grained text recognition,

Table 13: **Evaluation of existing LMMs on Chinese tasks of OCRBench v2' public data**. "LLM Size" indicates the number of parameters of the language model employed in each method.

| Method | LLM Size | Recognition | Extraction | Parsing | Understanding | Reasoning | Average |
|---|---|---|---|---|---|---|---|
| *Open-source LMMs* | | | | | | | |
| LLaVA-Next-8B [49] | 8B | 5.7 | 2.9 | 12.2 | 7.5 | 17.2 | 9.1 |
| LLaVA-OV-7B [50] | 8B | 14.8 | 15.7 | 13.7 | 16.0 | 28.7 | 17.8 |
| Monkey [51] | 8B | 4.6 | 11.2 | 8.4 | 21.5 | 20.0 | 13.1 |
| TextMonkey [7] | 8B | 23.5 | 14.8 | 8.4 | 19.9 | 12.2 | 15.8 |
| XComposer2-4KHD [140] | 7B | 16.7 | 18.8 | 12.1 | 27.5 | 2.3 | 15.5 |
| Molmo-7B [52] | 8B | 7.1 | 15.0 | 9.2 | 9.0 | 23.7 | 12.8 |
| Cambrian-1-8B [53] | 8B | 5.3 | 14.9 | 12.6 | 8.5 | 8.1 | 9.9 |
| Pixtral-12B [54] | 12B | 13.4 | 10.9 | 21.0 | 7.0 | 20.7 | 14.6 |
| EMU2-chat [141] | 37B | 2.3 | 0.5 | 8.5 | 1.0 | 7.3 | 3.9 |
| mPLUG-Owl3 [142] | 8B | 6.6 | 17.9 | 9.7 | 6.0 | 26.1 | 13.3 |
| CogVLM-chat [143] | 7B | 5.5 | 10.0 | 9.8 | 1.5 | 2.5 | 5.9 |
| Qwen-VL [4] | 8B | 7.2 | 5.3 | 10.7 | 11.5 | 11.2 | 9.2 |
| Qwen-VL-chat [4] | 8B | 9.5 | 8.2 | 9.3 | 11.0 | 21.1 | 11.8 |
| Qwen2-Vl-7B [66] | 7B | 51.3 | 51.4 | 21.6 | 52.5 | 37.5 | 42.9 |
| Qwen2.5-VL-7B [21] | 7B | **75.3** | 61.4 | **41.8** | 59.3 | 40.4 | 55.6 |
| InternVL2-8B [144] | 8B | 20.6 | 45.2 | 23.2 | 54.4 | 38.1 | 36.3 |
| InternVL2-26B [144] | 26B | 21.9 | 46.0 | 34.8 | 50.9 | 34.8 | 37.7 |
| InternVL2.5-8B [23] | 8B | 52.8 | 52.8 | 28.6 | 56.4 | 40.5 | 46.2 |
| InternVL2.5-26B [23] | 26B | 32.4 | 56.1 | 32.6 | 56.3 | 43.6 | 44.2 |
| InternVL3-8B [23] | 8B | 68.9 | 62.0 | 31.6 | 57.9 | 47.3 | 53.5 |
| InternVL3-14B [23] | 14B | 66.2 | **64.8** | 33.5 | **63.4** | **50.6** | **55.7** |
| Deepseek-VL-7B [145] | 7B | 8.0 | 13.3 | 15.7 | 5.5 | 18.5 | 12.2 |
| Deepseek-VL2-Small [55] | 16B | 60.9 | 50.6 | 28.3 | 53.0 | 20.5 | 42.7 |
| MiniCPM-V-2.6 [56] | 8B | 51.0 | 29.9 | 21.2 | 34.0 | 33.6 | 33.9 |
| MiniCPM-o-2.6 [56] | 7B | 53.0 | 49.4 | 27.1 | 43.5 | 32.7 | 41.1 |
| GLM-4V-9B [57] | 9B | 24.4 | 60.6 | 20.4 | 52.8 | 25.2 | 36.6 |
| VILA1.5-8B [146] | 8B | 5.4 | 8.8 | 8.5 | 3.0 | 15.5 | 8.2 |
| LLaVAR [30] | 13B | 2.3 | 1.7 | 8.9 | 0 | 2.5 | 3.1 |
| UReader [33] | 7B | 6.8 | 2.7 | 8.4 | 2.5 | 7.2 | 5.5 |
| DocOwl2 [147] | 7B | 4.2 | 10.3 | 8.6 | 4.0 | 9.6 | 7.3 |
| Yi-VL-6B [148] | 6B | 4.8 | 4.4 | 8.5 | 4.0 | 25.0 | 9.4 |
| Janus-1.3B [149] | 1.3B | 7.6 | 8.7 | 11.4 | 4.5 | 10.7 | 8.6 |
| Eagle-X5-7B [150] | 8B | 7.5 | 12.0 | 11.6 | 5.0 | 19.2 | 11.1 |
| Idefics3-8B [151] | 8B | 7.0 | 15.5 | 15.9 | 9.0 | 18.1 | 13.1 |
| Phi-4-MultiModal [152] | 5.6B | 51.5 | 32.3 | 12.1 | 34.4 | 23.0 | 30.7 |
| SAIL-VL-1.6-8B [153] | 8B | 31.2 | 40.0 | 23.9 | 42.3 | 35.0 | 34.5 |
| Kimi-VL-A3B-16B [154] | 16B | 57.2 | 54.7 | 31.5 | 52.5 | 31.4 | 45.5 |
| Ovis1.6-3B [58] | 3B | 11.5 | 23.7 | 22.8 | 28.8 | 18.9 | 21.1 |
| Ovis2-8B [58] | 7B | 72.2 | 50.8 | 37.7 | 47.9 | 37.4 | 49.2 |
| *Closed-source LMMs* | | | | | | | |
| GPT-4o [1] | - | 21.6 | 53.0 | 29.8 | 38.5 | 18.2 | 32.2 |
| GPT-4o-mini [59] | - | 13.1 | 38.9 | 27.2 | 28.8 | 16.9 | 25.0 |
| Gemini-Pro [60] | - | 52.5 | 47.3 | 30.9 | 51.5 | 33.4 | 43.1 |
| Claude3.5-sonnet [61] | - | 21.0 | 56.2 | 35.2 | 55.0 | 30.5 | 39.6 |
| Step-1V [62] | - | 56.7 | 41.1 | 37.6 | 38.3 | 39.2 | 42.6 |

text grounding, key information mapping, and document parsing. These figures show corresponding images and QA pairs for each of the 23 tasks. Fig. 17 to Fig. 18 provide examples of Chinese tasks, including key information extraction, text translation, formula recognition, reasoning VQA, cognition VQA, handwritten content extraction, document parsing, full-page OCR, and table parsing, along with their associated images and QA pairs.

## A.11 Samples for LMMs' Limitations

Fig. 19 to Fig. 21 provide examples corresponding to the findings discussed in Sec. 5.3 of the main text, which show error results of GPT-4o [1], Monkey [51], and Qwen2VL-8B on various tasks in *OCRBench v2*. These examples highlight the current limitations of LLMs on OCR tasks. For instance, LLMs exhibit poor recognition of less frequently encountered texts, struggle to accurately locate text in tasks involving text and coordinates, and demonstrate insufficient perception of text in complex layouts such as rotated texts. Additionally, their logical reasoning abilities are limited when addressing mathematical problems, and their analysis of complex elements in charts remains weak. These are the capabilities of LLMs in OCR tasks that require further improvement.

Table 14: **Evaluation of existing LMMs on English tasks of OCRBench v2's private data**.

| Method | Recognition | Referring | Spotting | Extraction | Parsing | Calculation | Understanding | Reasoning | Average |
|---|---|---|---|---|---|---|---|---|---|
| | | | Open-source LMMs | | | | | | |
| LLaVA-Next-8B [49] | 41.4 | 17.0 | 0 | 49.0 | 12.9 | 16.1 | 60.9 | 30.5 | 28.5 |
| LLaVA-OV-7B [50] | 45.4 | 18.5 | 0 | 60.0 | 15.5 | 32.0 | 59.0 | 39.3 | 33.7 |
| Monkey [51] | 31.5 | 0.1 | 0 | 34.4 | 26.3 | 17.7 | 61.4 | 22.4 | 24.2 |
| TextMonkey [7] | 39.8 | 1.6 | 0 | 27.6 | 24.8 | 10.2 | 62.3 | 21.2 | 23.4 |
| XComposer2-4KHD [140] | 39.5 | 12.0 | 0 | 69.7 | 26.0 | 20.2 | 68.2 | 35.8 | 33.9 |
| Molmo-7B [52] | 40.8 | 19.5 | 0 | 51.7 | 10.0 | 33.9 | 67.0 | 48.0 | 33.9 |
| Cambrian-1-8B [53] | 44.0 | 19.0 | 0 | 52.3 | 19.0 | 20.7 | 64.0 | 39.3 | 32.3 |
| Pixtral-12B [54] | 45.1 | 21.8 | 0 | 71.6 | 21.7 | 30.4 | 77.3 | 39.5 | 38.4 |
| EMU2-chat [141] | 34.3 | 0 | 0 | 20.4 | 21.3 | 20.3 | 47.1 | 18.3 | 20.2 |
| mPLUG-Owl3 [142] | 34.9 | 17.0 | 0 | 12.0 | 14.9 | 24.1 | 50.7 | 25.5 | 22.4 |
| CogVLM-chat [143] | 40.8 | 0 | 0 | 1.6 | 18.6 | 10.9 | 60.2 | 26.8 | 19.9 |
| Qwen-VL [4] | 35.9 | 4.2 | 0 | 38.7 | 28.5 | 13.8 | 60.1 | 16.9 | 24.8 |
| Qwen-VL-chat [4] | 34.1 | 12.6 | 0.1 | 42.6 | 19.5 | 18.4 | 58.3 | 20.3 | 25.7 |
| Qwen2-Vl-7B [66] | 47.0 | **42.0** | 1.5 | **90.2** | 13.7 | 36.4 | 71.1 | 36.6 | 42.3 |
| Qwen2.5-VL-7B [66] | 51.5 | 24.5 | <u>3.1</u> | 64.8 | 13.1 | 53.3 | <u>78.6</u> | 45.5 | 41.8 |
| InternVL2-8B [144] | 43.0 | 21.6 | 0 | 70.2 | 19.2 | 35.6 | 65.9 | 33.6 | 36.1 |
| InternVL2-26B [144] | 56.0 | 21.2 | 0 | 80.5 | 23.9 | 40.3 | 72.1 | 40.7 | 41.8 |
| InternVL2.5-8B [23] | 48.9 | 21.2 | 0 | 82.1 | 20.3 | 41.2 | 67.8 | 42.3 | 40.5 |
| InternVL2.5-26B [23] | 53.5 | 21.4 | 0 | 84.0 | 21.4 | 51.5 | 67.5 | 41.5 | 42.6 |
| InternVL3-8B [23] | 49.7 | 22.3 | 0.2 | 86.8 | 22.4 | 57.0 | 70.7 | 53.0 | 45.3 |
| InternVL3-14B [23] | 55.8 | 24.5 | 2.1 | 89.3 | 21.0 | <u>59.5</u> | 72.0 | 50.0 | 46.8 |
| Deepseek-VL-7B [145] | 33.5 | 13.7 | 0 | 19.1 | 11.7 | 24.8 | 60.5 | 32.5 | 24.5 |
| Deepseek-VL2-Small [55] | 56.6 | 23.7 | 0 | 86.4 | 18.9 | 30.6 | 72.2 | 39.5 | 41.0 |
| MiniCPM-V-2.6 [56] | 52.2 | 18.6 | 0.3 | 45.8 | 19.6 | 20.9 | 68.9 | 37.3 | 33.0 |
| MiniCPM-o-2.6 [56] | 54.1 | 24.7 | 0.3 | 74.4 | 17.6 | 39.2 | 75.7 | 47.0 | 41.6 |
| GLM-4v-9B [57] | 52.7 | 20.6 | 0 | 79.4 | 15.9 | 21.5 | 74.7 | 32.0 | 37.1 |
| VILA1.5-8B [146] | 36.0 | 14.5 | 0 | 26.0 | 17.4 | 20.3 | 44.7 | 27.0 | 23.2 |
| LLaVAR [30] | 13.8 | 0 | 0 | 8.3 | 15.2 | 4.4 | 42.4 | 15.0 | 12.4 |
| UReader [33] | 20.9 | 0 | 0 | 0 | 20.7 | 11.3 | 39.0 | 20.8 | 14.1 |
| DocOwl2 [147] | 25.4 | 7.5 | 0 | 47.1 | 26.2 | 8.3 | 52.8 | 19.5 | 23.4 |
| Yi-VL-6B [148] | 31.1 | 4.0 | 0 | 23.4 | 22.5 | 18.1 | 43.0 | 15.5 | 19.7 |
| Janus-1.3B [149] | 32.6 | 0 | 0 | 0.3 | 13.0 | 18.4 | 32.1 | 17.9 | 14.3 |
| Eagle-X5-7B [150] | 34.6 | 18.5 | 0 | 9.7 | 18.5 | 24.0 | 63.1 | 37.0 | 25.7 |
| Idefics3-8B [151] | 37.4 | 13.0 | 0 | 28.9 | 19.4 | 21.1 | 65.4 | 21.8 | 26.0 |
| Phi-4-MultiModal [152] | 58.4 | 19.0 | 0 | 53.5 | **38.7** | 28.7 | 66.8 | 39.8 | 38.1 |
| SAIL-VL-1.6-8B [153] | 56.7 | 24.1 | 2.2 | 79.3 | 22.8 | 45.4 | 69.2 | 45.3 | 43.1 |
| Kimi-VL-A3B-16B [154] | 49.1 | 13.5 | 0 | 28.8 | 21.9 | 37.6 | 69.4 | 36.2 | 32.1 |
| Ovis1.6-3B [58] | 48.5 | 19.5 | 0 | 69.2 | 20.7 | 22.1 | 74.6 | 49.5 | 38.0 |
| Ovis2-8B [58] | 54.2 | 20.9 | 0 | 83.6 | 24.2 | 54.7 | 74.1 | 57.3 | 46.1 |
| | | | Closed-source LMMs | | | | | | |
| GPT-4o [1] | <u>58.6</u> | 23.4 | 0 | 87.4 | 23.1 | 51.6 | 74.4 | **62.3** | <u>47.6</u> |
| GPT-4o-mini [59] | 55.3 | 21.8 | 0 | 85.4 | 20.6 | 45.2 | 75.5 | 49.0 | 44.1 |
| Gemini1.5-Pro [60] | **59.1** | <u>41.2</u> | **6.6** | <u>89.5</u> | 22.4 | 54.7 | **78.8** | <u>60.3</u> | **51.6** |
| Claude3.5-sonnet [61] | 52.9 | 24.9 | 2.5 | 86.9 | 23.8 | **61.4** | 74.4 | 53.0 | 47.5 |
| Step-1V [62] | 56.7 | 27.4 | 2.6 | 86.3 | <u>33.3</u> | 42.6 | 76.6 | 48.7 | 46.8 |

## A.12 Biases in Data Construction

Tab. 11 presents the scenario coverage statistics in our benchmark. The most frequent scenario accounts for 12.4% of the total samples. Among the 31 scenarios, 21 have more than 100 samples, which demonstrates the diversity of scene types in OCRBench v2.

In addition, we have manually verified all samples in our benchmark and did not identify any obvious regional or demographic biases.

## A.13 Broader Impacts

Our benchmark aims to enhance the evaluation of LMMs in text-oriented visual comprehension tasks. By establishing comprehensive benchmarks that reveal deficiencies in models' OCR capabilities, we provide insights for improving model performance. This advancement will elevate processing efficiency across scenarios such as document automation, assisted reading tools, and complex layout analysis, thereby benefiting applications in domains like healthcare and education. However, enhanced OCR functionality also introduces risks of misuse, including unauthorized extraction of sensitive information from images, surveillance-related applications, or generation of forged documents. To mitigate these risks, we restrict the use of this benchmark solely to research purposes and urge the community to prioritize privacy and fairness considerations in future model development.

Table 15: **Evaluation of existing LMMs on Chinese tasks of OCRBench v2's private data**.

| Method | LLM Size | Recognition | Extraction | Parsing | Understanding | Reasoning | Average |
|---|---|---|---|---|---|---|---|
| | | | Open-source LMMs | | | | |
| LLaVA-Next-8B [49] | 8B | 2.8 | 0.9 | 14.9 | 20.0 | 7.4 | 9.2 |
| LLaVA-OV-7B [50] | 8B | 5.4 | 13.6 | 20.3 | 34.0 | 13.6 | 17.4 |
| Monkey [51] | 8B | 1.5 | 28.4 | 29.1 | 40.0 | 8.3 | 21.5 |
| TextMonkey [7] | 8B | 10.5 | 15.2 | 30.2 | 44.0 | 7.6 | 21.5 |
| XComposer2-4KHD [140] | 7B | 12.9 | 38.6 | _37.5_ | 60.0 | 13.1 | 32.4 |
| Molmo-7B [52] | 8B | 3.4 | 29.8 | 6.6 | 24.0 | 11.1 | 15.0 |
| Cambrian-1-8B [53] | 8B | 2.4 | 19.8 | 26.7 | 36.0 | 7.6 | 18.5 |
| Pixtral-12B [54] | 12B | 6.2 | 22.3 | 11.4 | 26.0 | 14.0 | 16.0 |
| EMU2-chat [141] | 37B | 1.2 | 3.0 | 29.3 | 4.0 | 3.6 | 8.2 |
| mPLUG-Owl3 [142] | 8B | 1.6 | 27.4 | 27.3 | 16.0 | 10.0 | 16.5 |
| CogVLM-chat [143] | 7B | 2.4 | 16.2 | 22.5 | 20.0 | 3.1 | 12.8 |
| Qwen-VL [4] | 8B | 4.3 | 0 | 30.6 | 38.0 | 5.1 | 15.6 |
| Qwen-VL-chat [4] | 8B | 9.1 | 3.6 | 18.9 | 44.0 | 7.1 | 16.5 |
| Qwen2-Vl-7B [66] | 7B | 23.7 | 63.5 | 27.9 | 80.0 | 28.5 | 44.7 |
| Qwen2.5-VL-7B [66] | 8B | 24.4 | **78.9** | 33.1 | _82.0_ | 29.0 | 49.5 |
| InternVL2-8B [144] | 8B | 35.2 | 42.8 | 26.1 | 78.0 | 24.4 | 41.3 |
| InternVL2-26B [144] | 26B | 20.4 | 50.7 | 29.0 | 76.0 | 14.5 | 38.1 |
| InternVL2.5-8B [23] | 8B | 42.8 | 47.9 | 27.3 | 80.0 | 23.5 | 44.3 |
| InternVL2.5-26B [23] | 26B | 40.2 | 42.7 | 25.6 | 74.0 | 27.0 | 41.9 |
| InternVL3-8B [23] | 8B | 57.7 | 55.8 | 29.9 | 72.0 | 29.4 | 49.0 |
| InternVL3-14B [23] | 14B | 62.1 | 59.5 | 33.2 | 80.0 | 29.2 | 52.8 |
| Deepseek-VL-7B [145] | 7B | 3.2 | 14.7 | 10.7 | 30.0 | 9.8 | 13.7 |
| DeepSeek-VL2-Small [55] | 16B | 51.6 | 56.3 | 27.8 | 79.6 | 25.3 | 48.1 |
| MiniCPM-V-2.6 [56] | 8B | 53.1 | 53.2 | 32.8 | 76.0 | 23.4 | 47.7 |
| MiniCPM-o-2.6 [56] | 7B | 54.0 | 62.4 | 24.1 | 68.0 | 29.8 | 47.7 |
| GLM-4v-9B [57] | 9B | 60.6 | 65.2 | 32.4 | _82.0_ | 18.2 | 51.7 |
| VILA1.5-8B [146] | 8B | 1.4 | 9.1 | 22.2 | 16.0 | 6.4 | 11.0 |
| LLaVAR [30] | 13B | 2.2 | 2.0 | 27.1 | 10.0 | 1.9 | 8.6 |
| UReader [33] | 7B | 0.3 | 2.0 | 28.1 | 12.0 | 2.4 | 9.0 |
| DocOwl2 [147] | 7B | 1.0 | 17.8 | 29.4 | 20.0 | 3.9 | 14.4 |
| Yi-VL-6B [148] | 6B | 1.6 | 6.4 | 28.8 | 10.0 | 5.3 | 10.4 |
| Janus-1.3B [149] | 1.3B | 4.1 | 2.2 | 10.4 | 14.0 | 6.7 | 7.5 |
| Eagle-X5-7B [150] | 8B | 1.9 | 16.1 | 13.6 | 22.0 | 8.1 | 12.3 |
| Idefics3-8B [151] | 8B | 2.9 | 29.0 | 12.3 | 26.0 | 7.9 | 15.6 |
| Phi-4-MultiModal [152] | 5.6B | 30.5 | 40.5 | 42.7 | 56.0 | 16.9 | 37.3 |
| SAIL-VL-1.6-8B [153] | 8B | 35.8 | 41.5 | 35.7 | 76.0 | 23.9 | 42.6 |
| Kimi-VL-A3B-16B [154] | 16B | 54.0 | _71.1_ | 32.5 | **84.0** | 28.7 | 54.1 |
| Ovis1.6-3B [58] | 3B | 22.5 | 33.3 | 31.5 | 54.0 | 17.0 | 31.7 |
| Ovis2-8B [58] | 7B | 61.0 | 67.7 | **43.6** | _82.0_ | 25.6 | **56.0** |
| | | | Closed-source LMMs | | | | |
| GPT-4o [1] | - | 41.7 | 52.1 | 29.0 | 76.0 | 29.4 | 45.7 |
| GPT-4o-mini [59] | - | 20.0 | 53.6 | 27.9 | 66.0 | 19.6 | 37.4 |
| Gemini1.5-Pro [60] | - | **71.4** | 63.8 | 30.5 | _82.0_ | _29.9_ | _55.5_ |
| Claude3.5-sonnet [61] | - | 34.2 | 62.5 | 35.2 | 78.0 | **32.2** | 48.4 |
| Step-1V [62] | - | _65.2_ | 64.9 | 33.1 | 78.0 | 25.5 | 53.4 |

Table 16: Evaluation of InternVL2-8B with different resolution settings on the English tasks of OCRBench v2's public data.

| Method | Resolition | Recognition | Referring | Spotting | Extraction | Parsing | Calculation | Understanding | Reasoning | Average |
|---|---|---|---|---|---|---|---|---|---|---|
| | 448 | 47.3 | 19.1 | _0.1_ | 52.8 | **27.3** | 25.4 | 61.1 | 49.1 | 35.3 |
| InternVL2-8B [144] | 896 | _48.7_ | _23.0_ | **0.5** | **66.2** | _26.2_ | _25.9_ | _73.2_ | _51.9_ | _39.4_ |
| | dynamic | **49.9** | **23.1** | 0.5 | _65.2_ | 24.8 | **26.7** | **73.5** | **52.9** | **39.6** |

Table 17: Evaluation of Qwen2-VL-7B and Qwen2.5-7B with pre-provided OCR information on English tasks of OCRBench v2's public data.

| Method | Recognition | Referring | Spotting | Extraction | Parsing | Calculation | Understanding | Reasoning | Average |
|---|---|---|---|---|---|---|---|---|---|
| Qwen2-VL-7B [66] | **72.1** | _47.9_ | _17.5_ | **82.5** | _25.5_ | 25.4 | **78.4** | **61.5** | _51.4_ |
| Qwen2-VL-7B+OCR | _69.8_ | **50.4** | **20.1** | _79.1_ | **29.4** | _28.0_ | _77.7_ | 60.0 | **51.8** |
| Qwen2.5-8B+OCR | 28.6 | 13.8 | 0 | 45.9 | 24.2 | **31.3** | 61.1 | 40.5 | 30.7 |

Figure 10: Samples for each task.

## A.14 Limitations

One challenge we encountered is that LMMs sometimes produce responses that deviate from the given instructions, making it difficult to extract the desired answers. In future work, we plan to develop a more objective assessment framework to address this issue.

Another limitation arises when evaluating commercial LMMs, as some models occasionally refuse to answer certain questions due to safety filters or unclear content policies. This can lead to incomplete or biased performance assessments compared to open-source models that do not exhibit such behavior.

Figure 11: Samples for each task.

## Document Classification

What type of document is the image? Here is the options: letter, form, email, handwritten, advertisement, scientific report, scientific publication, specification, file folder, news article, budget, invoice, presentation, questionnaire, resume, memo. Directly output the results without any additional explanation.

Questionnaire

## Science QA

Based on the information from the provided image, answer the following question: Based on the Venn diagram, what do Natty Bumppo and Daniel Boone have in common? Here are the choices: 'Both fought in the French and Indian War.', 'Both were created by writer James Fenimore Cooper.'. Please output the answer directly.

Both fought in the French and Indian War.

Figure 12: Samples for each task.

## Chart Parsing

Convert the key information in the chart into a nested python dict.

```
{'title': 'Exploring Delicacies Around the World',
 'source': 'Source: Food and Travel Magazine',
 'x_title': 'None',
 'y_title': ['None', 'None'],
 'values': {
   'Cuisine Varieties': {
     'Africa': '150',
     'Asia': '180',
     'Europe': '200',
     'North America': '240',
     'South America': '300'
   },
   'Growth Rate': {
     'Africa': '10.0%',
     'Asia': '12.0%',
     'Europe': '15.0%',
     'North America': '18.0%',
     'South America': '20.0%'
   },
   'Local Preference': {
     'Africa': '25.0%',
     'Asia': '30.0%',
     'Europe': '35.0%',
     'North America': '40.0%',
     'South America': '45.0%'}}}
```

## Key Information Extraction

Find out the value of 'Serving Size', 'Calories/Energy per serving', 'Total Fat per serving', 'Total Fat % daily value/intake', 'Sodium per serving', 'Sodium % daily value/intake', 'Carbohydrate per serving', 'Carbohydrate % daily value/intake', 'Protein per serving' stated in the image. Extract the required information using the text in the image directly, and return the result in a dict with keys 'SS', 'CE-PS', 'TF-PS', 'TF-D', 'SO-PS', 'SO-D', 'CAR-PS', 'CAR-D', 'PRO-PS'.

```
{
  'SS': ['1/2 cup (122g)'],
  'CE-PS': ['90'],
  'TF-PS': ['0g'],
  'TF-D': ['0%'],
  'SO-PS': ['0mg'],
  'SO-D': ['0%'],
  'CAR-PS': ['21g'],
  'CAR-D': ['7%'],
  'PRO-PS': ['less than 1g']
}
```

Figure 13: Samples for each task.

## Full Page OCR



Read all the text in the image. Directly output the content and split the texts with space.

SURREALISM, ABSTRACTION\n& THE UNCONSCIOUS\nIn the 1920s, while a student at Yale, Rothko took\nphilosophy and psychology classes that led to his\ninterest in the unconscious and eventually to Sigmund\nFreud's The Interpretation of Dreams. Then, in the '40s,\nlike many American artists, he came under the influence\nof European Surrealist artists and writers living in New\nYork. Many artists also found inspiration in the Museum\nof Modern Art's influential 1936 exhibition, Fantastic Art,\nDada and Surrealism.\nBy the mid-1940s, Rothko loosened up his technique,\ninspired by the Surrealist method of automatism\n(\"automatic\" drawing or writing). Unlike some\nSurrealist artists, such as Salvador Dali, who \"pictured\"\nunconscious dreams in paintings, Surrealists using\nautomatism tried to access the unconscious by letting\nthe brush meander freely without planning or control.\nRothko experimented with the fluidity of watercolor and\nsoon realized he could achieve similar luminous effects\nin oil paint by diluting the pigment and applying color in\nthin washes, one on top of another. Rothko's imagery also\nchanged. Many works suggest paleontology and geology\nand evoke a vision of primordial life. Water seems to be\na primal element in which biomorphic shapes proliferate.\nSome compositions include stacked horizontal zones that\nmay stand for layers of the unconscious.

## Text Spotting



Spotting all the text in the image with word-level. Output the normalized coordinates of the left-top and right-bottom corners of the bounding box and the text content. The coordinates should be normalized ranging from 0 to 1000 by the image width and height.\nYour answer should be in the following format: [(x1, y1, x2, y2, text content), (x1, y1, x2, y2, text content)...] # The normalized coordinates and the content of the text in the image.

543, 770, 589, 794, 49-0223A,
545, 731, 580, 760, 502,
309, 594, 666, 641, YELLOWSTONE,
417, 160, 554, 198, TOUR

Figure 14: Samples for each task.

## Fine-grained Text Recognition



Recognize the text within the [192, 223, 332, 346] of the image. The coordinates have been normalized ranging from 0 to 1000 by the image width and height.

DIOS LE ABRE CAMINO\\n
AL HOMBRE\\n
QUE SABE A DONDE VA

## Text Grounding



Where is the region of the text 'COMNAM'? Output the normalized coordinates of the left-top and right-bottom corners of the bounding box. The coordinates should be normalized ranging from 0 to 1000 by the image width and height.
Your answer should be in the following format:
(x1, y1, x2, y2) # x1, y1, x2, y2 are the normalized coordinates of the bounding box.

[126,537,248,624]

## Key Information Mapping



According to the information in the image, please pair the corresponding keys and values below: Keys that need to be paired are 'Serving Size', 'Calories/Energy per 100g/ml', 'Carbohydrate per serving', 'Protein per 100g/ml', 'Total Fat per serving', 'Carbohydrate per 100g/ml', 'Total Fat per 100g/ml', 'Protein per serving'. Values that need to be paired are '0.8 g', '11.0 g', '200ml (1 cup)', '10.0 g', '1.6 g', '49 kcal(206 kJ)', '5.0 g', '5.5 g'.

{"Calories/Energy per 100g/ml": "49 kcal(206 kJ)"
  "Protein per serving": "10.0 g"
  "Protein per 100g/ml": "5.0 g"
  "Total Fat per serving": "1.6 g"
  "Total Fat per 100g/ml": "0.8 g"
  "Carbohydrate per serving": "11.0 g"
  "Carbohydrate per 100g/ml": "5.5 g"
  "Serving Size": "200ml (1 cup)"}

Figure 15: Samples for each task.

**Document Parsing**

Convert the privided document into markdown format.

We describe a winning strategy for Alice with $\Delta(G)$ colours in the $[B,A]$-edge colouring game played on $G$.\n\nThe only unsafe edges are the star edges of pending objects and the edge $vz$.\n\nAlice may arbitrarily number the pending objects $O_{1},O_{2},\ldots,O_{k+\ell}$ and performs basically the same pairing strategy as in the proof of Lemma 67 with only small extensions, as described in the following.\n\n* If Bob colours the matching edge of the pending object $O_{j}$, then, if this was the first such move and the edge $vz$ is still uncoloured, Alice colours $vz$ with the same colour (if possible, or a new colour otherwise); otherwise, Alice colours a star edge of the pending object $O_{j+1\mod{k+\ell}}$ with the same colour, if possible. If it is not possible, she uses a new colour for such a star edge.\n\n* If Bob colours the first star edge of the pending object $O_{j}$ and there is still a pending object with only uncoloured star edges, then Alice colours the matching edge of the pending object $O_{j-1\mod{k+\ell}}$ with the same colour. If the matching edge is already coloured, then Alice misses her turn.\n\n* If Bob colours the first star edge of the pending object $O_{j}$ and there is no pending object with only uncoloured star edges left, then Alice colours $vz$ with a new colour (if $vz$ is still uncoloured) or misses her turn (if $vz$ is already coloured).\n\n* If Bob colours the edge $vz$, an edge $vx_{j}$ or the second star edge of the pending object (a triangle) $O_{j}$, then Alice misses her turn.\n* If Bob colours an edge $zu_{i}$, then Alice colours $vz$ (if $vz$ is still uncoloured) or misses her turn (otherwise).\n\nThis strategy has the same properties as the strategy for the single galaxy in the proof of Lemma 67, and, in addition, it guarantees that the edge $vz$ is coloured before it is in danger to be infeasible for any colour. \n\n### Permitted for game $[A,a]$\n\n**Definition 69** (full tree).: Let $n,m_{1},m_{2}\in\mathbb{N}$. An _$(n,m_{1},m_{2})$-full tree_ is based on a path $P_{3}$, where there are $m_{1}$ (respectively, $n$), $m_{2}$) leafs attached its three vertices, i.e., the graph has the vertex set\n\n$$\{w_{1},v,w_{2}\}\cup\{x_{i}\mid 1\leq i\leq m_{1}\}\cup\{y_{j}\mid 1\leq j\leq n\}\cup\{z_{i}\mid 1\leq i\leq m_{2}\}$$\n\nand the edge set\n\n$$\{w_{1}v,vw_{2}\}\cup\{w_{1}x_{i}\mid 1\leq i\leq m_{1}\}\cup\{vy_{j}\mid 1 \leq j\leq n\}\cup\{w_{2}z_{i}\mid 1\leq i\leq m_{2}\}.$$\n\nA _full tree_ is an $(n,m_{1},m_{2})$-full tree for some $n,m_{1},m_{2}\in\mathbb{N}$."

Figure 16: Samples for each task.

Figure 17: Samples for each task.

## Document Parsing

**[User]** Parse the document image in Markdown format.

**[Assistant]** 西安石油大学硕士学位论文    \n图4-5 2ED020I12FA型号IGBT驱动器的实物图及内部原理图    \na)实物图 b)内部原理图    \n表4-5 2ED020I12FA型号IGBT驱动器部分引脚功能    \n<table><tr><td>引脚序号</td><td>名称</td></tr></td>功能</td></tr><tr><td>2</td><td>INHS+</td><td>非反相驱动器输入高端</td></tr><tr><td>3</td><td>INHS-</td><td>反向驱动器输入高端</td></tr><tr><td>7</td><td>VCC1HS</td><td>正电源输入高端</td></tr><tr><td>12</td><td>INLS+</td><td>非反相驱动器输入低端</td></tr><tr><td>13</td><td>INLS-</td><td>反相驱动器输入低端</td></tr></table> \n\n44  \n(C)1994-2020 China Academic Journal Electronic Publishing House.All rights reserved.http://www.cnki.net

## Full Page OCR

**[User]** Read all the text in the image .

**[Assistant]** 着力打造上海国际旅游度假区,
提升上海旅游业的国际竞争力和影响力,
推动世界著名旅游城市建设。

## Table Parsing

**[User]** Parse the HTML-formatted structural information of all the tables in the image.

**[Assistant]** <table><tr><td> 星期 </td><td> 一 </td><td> 二</td><td> 三</td><td> 四</td><td> 五</td></tr> <tr><td> 每股涨跌</td><td> +4</td><td> +4.5</td><td> -1</td><td> -2.5</td><td> -6</td> </tr></table>

Figure 18: Samples for each task.

Figure 19: **Samples for LMM'S Limitations**. The portion of LLM's response marked in red is incorrect content, and the content in the red dashed box is missing information.

Figure 20: **Samples for LMM'S Limitations**. The portion of LLM's response marked in red is incorrect content, and the content in the red dashed box is missing information.

Figure 21: **Samples for LMM'S Limitations**. The portion of LLM's response marked in red is incorrect content, and the content in the red dashed box is missing information.