# TRIG-BENCH:
# A BENCHMARK FOR TEXT-RICH IMAGE GROUNDING

**Ming Li**[1]**, Ruiyi Zhang**[2]**, Jian Chen**[3]**, Tianyi Zhou**[1]
[1]University of Maryland     [2]Adobe Research     [3]University at Buffalo
minglii@umd.edu, ruizhang@adobe.com
https://github.com/MingLiiii/TRIG-Bench

## ABSTRACT

Despite the existing evolution of Multimodal Large Language Models (MLLMs), a non-neglectable limitation remains in their struggle with visual text grounding, especially in text-rich images of documents. Document images, such as scanned forms and infographics, highlight critical challenges due to their complex layouts and textual content. However, current benchmarks do not fully address these challenges, as they mostly focus on visual grounding on natural images, rather than text-rich document images. Thus, to bridge this gap, we introduce **TRIG**, a novel task with a newly designed instruction dataset for benchmarking and improving the **T**ext-**R**ich **I**mage **G**rounding capabilities of MLLMs in document question-answering. Specifically, we propose an OCR-LLM-human interaction pipeline to create 800 manually annotated question-answer pairs as a benchmark **TRIG-Bench** based on four diverse datasets. A comprehensive evaluation of various MLLMs on our proposed benchmark exposes substantial limitations in their grounding capability on text-rich images.

## 1 INTRODUCTION

Despite the remarkable advancements in LLMs (Chang et al., 2024; Hadi et al., 2023; Xu et al., 2024; Thirunavukarasu et al., 2023) and MLLMs (Zhang et al., 2024a; Ghosh et al., 2024), the trustworthiness of their generated outputs remains a critical concern (Liu et al., 2024d; Sun et al., 2024a). While these models can produce fluent and coherent responses, they often lack grounding capability, which can lead to potential hallucinations (Ji et al., 2023; Rawte et al., 2023; Liu et al., 2024a). Grounding capability is defined as the model's ability to accurately localize relevant regions in the visual content based on the provided semantic description (Nagaraja et al., 2016; Luo & Shakhnarovich, 2017; Yu et al., 2017; Kamath et al., 2021; You et al., 2024). This capability is essential for ensuring that the responses are accurate and verifiable. By providing grounding information, these models enable real-world users to easily verify the correctness of the generated responses, thereby mitigating the uncertainty associated with LLM outputs. This grounding serves as a crucial bridge for enhancing the interaction between humans and LLMs, fostering greater trust and transparency in AI-generated content.
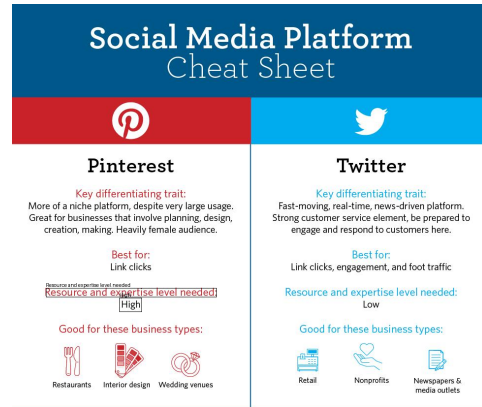


Figure 1: An example from Info-grahicsVQA. **Question:** What is the resource and expertise level needed for Pinterest? **Answer:** high. The LLM is expected to generate the answer together with the corresponding grounded bounding boxes that can support its answer, which requires deeper spatial understanding and reasoning, and sometimes instruction-following abilities.

---

This work was done at the University of Maryland.

Existing grounding efforts in MLLMs mainly focus on natural images, where the task involves associating textual descriptions with corresponding visual elements, such as objects or scenes(Plummer et al., 2015; Chen et al., 2023b; Zhang et al., 2023a; 2024c; Peng et al., 2024; You et al., 2024). However, there is a significant gap in the literature when it comes to grounding in text-rich document images. Document images, such as scanned forms, charts, and complex posters, present unique challenges that differ markedly from those found in natural images (Mathew et al., 2021; Masry et al., 2022; Mathew et al., 2022; Zhang et al., 2024b). They often feature a mix of textual and graphical elements and require precise localization and understanding of content, especially textual content. An example is shown in Figure 1, for the given text-rich document image, we expect LLMs to not only generate the answer alone but also to provide the corresponding grounded bounding boxes that can support its answer, which requires deeper spatial understanding and reasoning capabilities, and sometimes better instruction-following abilities to correctly provide the grounding information in the desired formats.

Despite its importance, there is no established benchmark specifically designed to evaluate the grounding capabilities of MLLMs on Text-Rich Document Question-Answering tasks. The absence of such a benchmark and training data limits the ability to systematically assess and improve the performance of different models in this domain. Thus to bridge this gap, we introduce **TRIG**, a novel **T**ext-**R**ich **I**mage **G**rounding task with instruction set for document QA grounding, along with its corresponding benchmark notated as **TRIG-Bench**. TRIG-Bench consists of 800 question-answer pairs manually collected from DocVQA (Mathew et al., 2021), ChartQA (Masry et al., 2022), InfographicsVQA (Mathew et al., 2022), and TRINS datasets (Zhang et al., 2024b), along with human-inspected ground-truth bounding boxes that support the answer to the corresponding question. They provide a standardized framework for evaluating the ability of MLLMs to accurately ground their responses to document-related visual questions.

For text-rich document images, we mainly focus on visual texts on them as the main grounding target. Considering the supreme performance of the modern OCR models (Chen et al., 2021; Subramani et al., 2021; Du et al., 2022; Li et al., 2023c;b; Wei et al., 2024; Sun et al., 2024b) and the promising reasoning ability of current LLMs like GPT4o, an OCR-LLM-Human interaction pipeline is proposed for the benchmark construction. Specifically, for every given VQA pair, we first utilize PaddleOCR to detect and recognize all the texts. Given all these OCR results, LLMs are prompted to judge which bounding boxes support the answer to the corresponding question, followed by another LLM evaluating the correctness of the grounded bounding boxes chosen. The resulting data has already been of good quality as it is generated and verified by the powerful GPT4o model. Then, further human participants will manually inspect the generated results and select the correct ones to form our benchmark data. Our main contributions:

- We introduce a novel benchmark, **TRIG-Bench**, specifically designed for the challenging text-rich document image grounding task. This benchmark is the first of its kind and provides a standardized framework for evaluating MLLMs in this domain, filling a critical gap in existing research.
- We conduct a comprehensive evaluation of a range of existing MLLMs on our new benchmark. Our analysis provides a deeper understanding of the limitations and constraints of current MLLMs when applied to this new and complex task.
- Our findings reveal a significant issue: Although most of the current MLLMs perform well on well-defined tasks, they lack the capability to follow customized and complex tasks that require a deep understanding and reasoning capabilities.

## 2 RELATED WORKS

Grounding (Harnad, 1990) is essential for effective human communication with machines. Currently, grounding datasets have been utilized for vision-language pre-training and improvement for both object-level recognition (Li et al., 2022) and language acquisition (Ma et al., 2023). Recent studies propose to integrate text and grounding regions into token sequences (Yang et al., 2022; Lu et al., 2022; Wang et al., 2022) within language modeling frameworks. Building on such approaches, researchers have developed a series of grounded MLLMs, including GPT4ROI (Zhang et al., 2023a), Kosmos-2 (Peng et al., 2023), Shikra (Chen et al., 2023b), PVIT (Chen et al., 2023a), BuboGPT (Zhao et al., 2023), Qwen-VL (Bai et al., 2023b), and Ferret (You et al., 2023). Despite

their impressive performance, these models mainly focus on the grounding of natural image objects, and the grounding of text-rich document images is still under-explored. For the two most related works which also mentioned text grounding on document images, P²G (Chen et al., 2024) includes the OCR model into the QA pipeline as an information amplifier; TG-Doc (Wang et al., 2023b) utilizes grounding capability to improve models QA capabilities. Both of these two works treat the text grounding process as an intermediate step and do not conduct any evaluation of the grounding capability. On the contrary, we entirely focus on the grounding capability and first propose the benchmark for it.

## 3 DATA CONSTRUCTION

In this section, we focus on the construction of our benchmark and training dataset, including the Data Source, the Construction Pipeline, and Statistics.

### 3.1 DATA SOURCE

To ensure the diversity of our benchmark data, four existing document datasets are chosen as our data sources, covering a wide range of document image types, including **DocVQA** (Mathew et al., 2021), **ChartQA** (Masry et al., 2022), **InfographicVQA** (Mathew et al., 2022), and **TRINS** (Zhang et al., 2024b).

### 3.2 CONSTRUCTION PIPELINE

Our OCR-LLM-Human interactive pipeline for training data and benchmark data generation is shown in Figure 2.

**Step 1 Preprocessing:** PaddleOCR[1] is utilized to obtain the initial OCR information for its simplicity and promising performance.

**Step 2 Generation:** After obtaining the OCR information, a critical issue is how to transmit the OCR information to the LLMs. The common method in the literature is wrapping all the obtained OCR information to the prompt (Zhang et al., 2023b; Wang et al., 2023b), however, we observe a severe misalignment between the text information in the prompt and the visual information in the image, leading to unpromising results. Thus to further align the visual and text information provided in different genres, we innovatively assign every OCR bounding box an identical index, draw every bounding box with its index on the original images, and simultaneously provide all the OCR information with the index in the prompt. By utilizing this strategy, the LLM is able to better align each visual element with its detailed information. Along with the original question and ground truth answer, the LLM is prompted to select the bounding boxes that can support the given answer to the corresponding question.
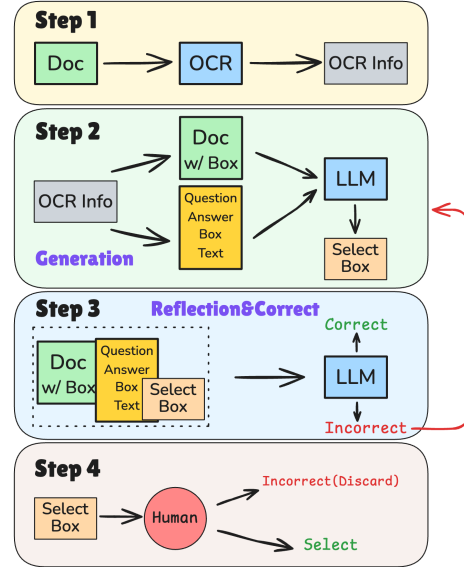


Figure 2: Main Constriction Pipeline. The pipeline contains 4 steps: Preprocessing, Generation, Correction, and Human Evaluation. The benchmark data will go through all of these 4 steps and the training data will go through the previous 3 steps.

**Step 3 Correction:** Then another Reflection & Rectification module (Pan et al., 2023; Huang et al., 2023; Li et al., 2023a; 2024b) is introduced to examine the correctness of selected bounding boxes. In addition to all the previous information used in Step 2, the previously generated bounding box indices are also provided to the LLM, prompting it to judge if the selected grounded bounding boxes can adequately lead to the given answer. If the judging result is "correct", this sample will be kept

---

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

if it is from the training set, or sent to the next step if it is used for building benchmark. Otherwise, it will be sent back to the previous step until reaching the maximum round and being discarded. The resulting data after the previous 3 steps has already been of good quality as it is generated and verified by the powerful LLM.

**Step 4 Human Evaluation:** The sample sent to step 4 has been evaluated and rectified by LLM, however, its quality might still have discrepancies with the requirements of being a benchmark. Thus to ensure the correctness of our benchmark data, further human inspection is conducted. Specifically, two human participants are invited for the human evaluation. Only those samples that are agreed to be correct for both of the participants can be kept in our benchmark, otherwise will be discarded. Finally, for each dataset, 200 QA pairs are manually selected as the ground truth of this benchmark.

### 3.3 DATA STATISTICS

The benchmark data statistics are presented in Table 1. "Avg Question Len", "Avg Answer Len", and "Avg OCR Text Len" calculate the average question, answer, and OCR text word counts. "Avg OCR Box #", "Avg GT Box", and "Avg GT Box Ratio" calculate the number of bonding boxes provided by OCR models, selected as the ground truth and their ratio. "Avg OCR Area (%)", "Avg GT Area (%)" and "Avg GT Area Ratio" calculate the area of bonding boxes provided by OCR models, selected as the ground truth and their ratio. From these statistics, the characteristics of each dataset can be illustrated, which showcases the variety of our data components and further provides a clear understanding of evaluation results.

|  | Chart | Doc | Info | Trins | Total |
|---|---|---|---|---|---|
| Total Question # | 200 | 200 | 200 | 200 | 800 |
| Total Image # | 171 | 190 | 199 | 199 | 759 |
| Avg Question Len | 11.04 | 9.66 | 12.42 | 11.69 | 11.20 |
| Avg Answer Len | 1.39 | 2.94 | 1.93 | 19.48 | 6.44 |
| Avg OCR Text Len | 2.7 | 4.73 | 3.02 | 3.00 | 3.36 |
| Avg OCR Box # | 26.15 | 53.30 | 102.26 | 8.94 | 47.66 |
| Avg GT Box # | 2.67 | 1.73 | 2.78 | 2.45 | 2.41 |
| Avg GT Box Ratio | 11.83% | 4.88% | 3.72% | 37.74% | 14.54% |
| Avg OCR Area (%) | 12.56 | 19.75 | 18.29 | 15.82 | 16.61 |
| Avg GT Area (%) | 0.80 | 0.79 | 0.53 | 6.09 | 2.05 |
| Avg GT Area Ratio | 8.59% | 5.10% | 2.89% | 42.69% | 14.82% |

Table 1: **Benchmark Data Statistics.** Total Question # represents the unique question number from each dataset. Total Image # represents the Unique image number. Other statistics are averaged on each dataset.

## 4 EXPERIMENTAL RESULTS

For the instruction-based method, we utilize LLaVA-v1.5-Vicuna-13B (Liu et al., 2023a) as our base model to be fine-tuned. For the embedding-based method, we utilize PaliGemma-3B (Beyer et al., 2024) as our base model to be fine-tuned. Detailed training configurations, data statistics, and examples can be found in the supplementary material. All experiments are performed on Nvidia A100 GPUs.

### 4.1 MAIN RESULTS

The overall evaluation results across various MLLMs are presented in Table 2 (Evaluation Setting 1), Table 3 (Evaluation Setting 2). The MLLMs we evaluated are listed as follows: LLaVA-v1.6-Vicuna-13B, LLaVA-v1.6-Vicuna-7B (Liu et al., 2023a;b; 2024b), Phi3-V(Team, 2024), DeepSeek-VL-7B-chat (Lu et al., 2024), Idefics2-8B (Laurençon et al., 2023; 2024), Qwen-VL (Bai et al., 2023a), CogVLM2-Llama3-19B (Wang et al., 2023a), InternLM-XComposer2-VL-7B (Dong et al., 2024a), InternLM-XComposer2-4KHD-7B (Dong et al., 2024b), Monkey-Chat (Li et al., 2024d), MiniCPM-Llama3-V 2.5 (Yao et al., 2024), and GPT series.

As shown in Table 2, all the existing models, including open-source models and the powerful GPT4o, perform not well for the **OCR-free Grounding** setting, in which they are required to gen-

| Testsets | Chart | Doc | Info | Trins | Avg |
|---|---|---|---|---|---|
| Metrics (%) | IoU | IoU | IoU | IoU | Avg |
| LLaVA-v1.6-Vicuna-13B | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| LLaVA-v1.6-Vicuna-7B | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| Phi3-V | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| DeepSeek-VL-7B-chat | 0.07 | 0.00 | 0.02 | 0.00 | 0.02 |
| Idefics2-8B | 0.21 | 0.01 | 0.00 | 0.00 | 0.06 |
| Qwen-VL | 0.43 | 0.06 | 0.18 | 0.23 | 0.22 |
| CogVLM2-Llama3-19B | 0.19 | 0.01 | 0.16 | 0.66 | 0.25 |
| InternLM-XComposer2-VL-7B | 0.15 | 0.20 | 0.13 | 0.57 | 0.26 |
| Monkey-Chat | 0.77 | 0.19 | 0.15 | 0.45 | 0.39 |
| InternLM-XComposer2-4KHD-7B | 1.04 | 0.10 | 0.90 | 0.14 | 0.55 |
| MiniCPM-Llama3-V 2.5 | 0.44 | 1.40 | 0.65 | 4.96 | 1.86 |
| GPT-4o | 3.90 | 1.79 | 1.60 | 13.73 | 5.26 |

Table 2: **OCR-free Grounding.** Chart, Doc, Info, and Trins represent evaluation results on ChatQA, DocVQA, InfographicsVQA, and TRINS datasets, respectively. IoU represents the pixel-level IoU scores. Avg represents the average IoU score on the 4 datasets and the ordering of each model is decided by this average score.

| Testsets | ChartQA | | | | DocVQA | | | | InfographicsVQA | | | | TRINS | | | | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics (%) | IoU | P | R | F1 | IoU | P | R | F1 | IoU | P | R | F1 | IoU | P | R | F1 | Avg |
| LLaVA-v1.6-Vicuna-13B | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| LLaVA-v1.6-Vicuna-7B | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| Idefics2-8b | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.17 | 1.17 | 3.50 | 1.69 | 0.47 |
| DeepSeek-VL-7B-chat | 0.84 | 1.35 | 2.25 | 1.23 | 0.19 | 0.19 | 1.50 | 0.34 | 0.00 | 0.00 | 0.00 | 0.00 | 1.60 | 1.60 | 2.00 | 1.67 | 0.92 |
| InternLM-XComposer2-4KHD-7B | 1.00 | 2.00 | 1.00 | 1.25 | 0.25 | 0.50 | 0.50 | 0.33 | 0.00 | 0.00 | 0.00 | 0.00 | 3.67 | 6.80 | 3.79 | 4.65 | 1.59 |
| Monkey-Chat | 3.58 | 5.38 | 9.91 | 5.05 | 0.94 | 1.19 | 1.75 | 1.15 | 0.34 | 0.30 | 2.00 | 0.51 | 0.00 | 0.00 | 0.00 | 0.00 | 2.01 |
| Phi3-V | 2.54 | 3.14 | 5.46 | 3.25 | 0.97 | 1.22 | 2.00 | 1.21 | 0.40 | 0.35 | 2.08 | 0.59 | 3.81 | 4.31 | 5.00 | 4.28 | 2.54 |
| CogVLM2-Llama3-19B | 1.65 | 2.30 | 3.68 | 2.02 | 1.56 | 1.87 | 3.58 | 1.89 | 0.42 | 0.84 | 1.75 | 0.71 | 6.76 | 7.93 | 8.33 | 7.53 | 3.30 |
| Qwen-VL | 4.36 | 4.33 | 28.00 | 7.27 | 1.07 | 1.06 | 8.25 | 1.83 | 0.70 | 0.71 | 4.50 | 1.20 | 1.51 | 1.83 | 3.32 | 2.08 | 4.50 |
| InternLM-XComposer2-VL-7B | 8.10 | 15.31 | 9.03 | 10.49 | 7.36 | 11.28 | 8.08 | 8.73 | 2.32 | 6.39 | 2.85 | 3.39 | 15.82 | 18.12 | 17.32 | 16.98 | 10.10 |
| MiniCPM-Llama3-V 2.5 | 7.40 | 12.50 | 8.40 | 9.24 | 15.78 | 19.22 | 18.42 | 17.48 | 5.71 | 9.87 | 7.34 | 7.14 | 54.06 | 62.06 | 57.89 | 57.95 | 23.15 |
| GPT-4o | 83.80 | 88.80 | 89.24 | 87.47 | 82.14 | 87.14 | 89.50 | 86.16 | 68.19 | 79.57 | 78.81 | 75.82 | 89.08 | 96.06 | 91.53 | 92.16 | 85.34 |
| GPT-3.5-turbo (Without Image) | 8.81 | 14.17 | 9.48 | 10.64 | 32.05 | 40.92 | 32.50 | 35.00 | 12.66 | 20.45 | 14.39 | 15.75 | 15.81 | 18.75 | 15.80 | 16.62 | 19.58 |
| GPT-4 (Without Image) | 51.58 | 57.29 | 53.75 | 54.34 | 52.18 | 62.79 | 53.61 | 56.28 | 47.31 | 57.51 | 54.34 | 53.51 | 69.83 | 76.16 | 71.05 | 72.39 | 58.59 |
| GPT-4o (Without Image) | 59.50 | 67.13 | 62.84 | 63.29 | 77.83 | 83.41 | 80.71 | 80.80 | 63.34 | 72.64 | 69.03 | 68.88 | 71.05 | 80.01 | 72.54 | 74.38 | 71.69 |

Table 3: **OCR-based Grounding.** IoU, P, R, F1 represent bounding-box-level IoU score, precision, recall and F1 score. Avg represents the average score on all datasets and evaluation metrics and the ordering is decided by this score.

erate supportive bounding boxes without any additional bounding box information. From the evaluation setting 1 results, we can conclude that the capability of existing models to generate grounded bounding boxes from scratch is limited and thus the training specifically for this setting is required.

As shown in Table 3, for the **OCR-based Grounding** setting, most of the models can reach a much better performance as the OCR information has been provided. However, the performances of open-source models still have a large gap with GPT4o.

## 4.2 ANALYSIS

> **Finding 1.** All existing MLLMs are not good at generating grounded bounding boxes from scratch.

As shown in Table 2, even the most powerful GPT4o can only gain an average IoU score of $5.28\%$, consisting of $13.73\%$ on TRINS and approximately under $4.0\%$ on other datasets. The relatively higher performance on TRINS is due to its special characteristic that the images in TRINS contain the least number of OCR objects while occupying the most area, making it easier to generate intersected bounding boxes[2]. However, when it comes to common information-intense documents, the performances drop dramatically as the ground truth bounding boxes become much smaller.

Compared with GPT4o, the IoU values from other open-source models are mostly under $1.0\%$, which is negligible. GPT4o is able to follow the instructions and generate bounding boxes, though its relatively weak spatial understanding makes the generation not precise enough. However, most of the other open-source models are not able to either understand our instructions or generate reasonable bounding boxes, especially for those MLLMs that get near-zero IoU values. Even if

---

[2]Examples can be found in the supplementary material for better understanding.

we provide further instructions requiring them to generate bounding boxes that can support their answer, they are not able to understand the instructions nor to follow them.

These results reveal a critical issue that current MLLMs have a relatively weak spatial understanding and are not capable of generating grounded boxes that support their answer from scratch, which makes it a potential future direction.

> **Finding 2.** Most existing open-source MLLMs are not able to follow customized complex instructions.

Table 3 represents the evaluation setting where the OCR information, including bounding boxes and texts, is wrapped into the input to MLLMs, thus making the whole process a much simpler bounding box selection process. The performance of GPT4o reaches an astonishingly high value of $85.34\%$ on average compared with the $5.26\%$ on evaluation setting 1, indicating that generating grounded bounding boxes from scratch is hard for GPT4o. In this setting, even the text-only models can achieve a reasonably high performance due to the detailed information provided by OCR models while the performances of most of the existing open-source MLLMs still kept low, making it impossible for practical usage.

By careful inspection, we observe that these low performances on existing open-source MLLMs are mainly caused by their inability to follow the given instructions: most existing open-source MLLMs will directly generate corresponding answers to the question and ignore the instruction of selecting supporting grounded bounding boxes, potentially due to the overfitting on the format of the training data. To further quantitatively analyze this issue, we introduce another value, **the instruction-following rate**, defined by the proportion of testing samples for which the MLLMs are able to generate at least one bounding box required in the additional instruction. This value does not measure the correctness of MLLM-generated bounding boxes but only the existence of them. It directly represents MLLMs' instruction-follow ability for this task, i.e. generate at least one bounding box regardless of the correctness.

From the results of the instruction-following rates, we can see that GPT4o reaches an average instruction-following rate of $98\%$. Even if for the OCR-free Grounding setting, in which it can only reach an average IoU score of $5.26\%$, it still achieves the instruction-following rate of $96\%$. The comparison between the instruction-following rate and IoU score reveals that (i) GPT4o has a strong capability in following customized complex instructions; (ii) The low IoU is caused by GPT4o's inability of spatial understanding. On the contrary, the instruction-following rates on existing open-source MLLMs are mostly less than $30.0\%$ or even $10.0\%$, which means under most circumstances, they do not understand our instruction and do not generate any bounding boxes, representing a "not-even-wrong" situation. Since the average instruction-following rates are much higher than the grounding metrics, the gaps between them indicate the MLLMs' inability to ground.

Our settings reflect both MLLMs' instruction-following and grounding ability. Most of the existing open-source models are not able to follow the relatively complex instructions in this task, and even if they correctly follow the instructions, most of them are not able to find the correct grounded bounding boxes that support their answers. Comparing these results with their outstanding performances on standard QA tasks, we reveal this critical issue that most existing open-source MLLMs are potentially overfitted to the standardized tasks, and they still lack instruction-following abilities.

## 5    CONCLUSION

In this paper, we introduce **TRIG**, a novel task designed to evaluate and enhance the visual text grounding capabilities of MLLMs in text-rich document images. We first build **TRIG-Bench**, a comprehensive benchmark. Our evaluation reveals significant limitations in the visual text grounding capabilities of existing MLLMs, particularly in their ability to handle complex document layouts and textual content.

## REFERENCES

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023a.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. 2023b.

Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, Thomas Unterthiner, Daniel Keysers, Skanda Koppula, Fangyu Liu, Adam Grycner, Alexey Gritsenko, Neil Houlsby, Manoj Kumar, Keran Rong, Julian Eisenschlos, Rishabh Kabra, Matthias Bauer, Matko Bošnjak, Xi Chen, Matthias Minderer, Paul Voigtlaender, Ioana Bica, Ivana Balazevic, Joan Puigcerver, Pinelopi Papalampidi, Olivier Henaff, Xi Xiong, Radu Soricut, Jeremiah Harmsen, and Xiaohua Zhai. Paligemma: A versatile 3b vlm for transfer, 2024. URL https://arxiv.org/abs/2407.07726.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45, 2024.

Chi Chen, Ruoyu Qin, Fuwen Luo, Xiaoyue Mi, Peng Li, Maosong Sun, and Yang Liu. Position-enhanced visual instruction tuning for multimodal large language models. *arXiv preprint arXiv:2308.13437*, 2023a.

Jiaxing Chen, Yuxuan Liu, Dehu Li, Xiang An, Weimo Deng, Ziyong Feng, Yongle Zhao, and Yin Xie. Plug-and-play grounding of reasoning in multimodal large language models, 2024. URL https://arxiv.org/abs/2403.19322.

Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm's referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023b.

Xiaoxue Chen, Lianwen Jin, Yuanzhi Zhu, Canjie Luo, and Tianwei Wang. Text recognition in the wild: A survey. *ACM Comput. Surv.*, 54(2), March 2021. ISSN 0360-0300. doi: 10.1145/3440756. URL https://doi.org/10.1145/3440756.

Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang Zhang, Haodong Duan, Maosong Cao, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Xinyue Zhang, Wei Li, Jingwen Li, Kai Chen, Conghui He, Xingcheng Zhang, Yu Qiao, Dahua Lin, and Jiaqi Wang. Internlm-xcomposer2: Mastering free-form text-image composition and comprehension in vision-language large model. *arXiv preprint arXiv:2401.16420*, 2024a.

Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Songyang Zhang, Haodong Duan, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Zhe Chen, Xinyue Zhang, Wei Li, Jingwen Li, Wenhai Wang, Kai Chen, Conghui He, Xingcheng Zhang, Jifeng Dai, Yu Qiao, Dahua Lin, and Jiaqi Wang. Internlm-xcomposer2-4khd: A pioneering large vision-language model handling resolutions from 336 pixels to 4k hd. *arXiv preprint arXiv:2404.06512*, 2024b.

Yongkun Du, Zhineng Chen, Caiyan Jia, Xiaoting Yin, Tianlun Zheng, Chenxia Li, Yuning Du, and Yu-Gang Jiang. Svtr: Scene text recognition with a single visual model. 2022. URL https://arxiv.org/abs/2205.00159.

Akash Ghosh, Arkadeep Acharya, Sriparna Saha, Vinija Jain, and Aman Chadha. Exploring the frontier of vision-language models: A survey of current methodologies and future directions, 2024. URL https://arxiv.org/abs/2404.07214.

Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili, et al. Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects. *Authorea Preprints*, 2023.

Stevan Harnad. The symbol grounding problem. *Physica D: Nonlinear Phenomena*, 42(1-3):335–346, 1990.

Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large language models can self-improve. In Houda Bouamor, Juan Pino, and Kalika Bali

(eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 1051–1068, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.67. URL `https://aclanthology.org/2023.emnlp-main.67/`.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, 2023.

Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1780–1790, 2021.

Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. Obelics: An open web-scale filtered dataset of interleaved image-text documents, 2023.

Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models?, 2024.

Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10965–10975, 2022.

Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, and Tianyi Zhou. Reflection-tuning: Recycling data for better instruction-tuning. In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*, 2023a. URL `https://openreview.net/forum?id=xaqoZZqkPU`.

Ming Li, Bin Fu, Han Chen, Junjun He, and Yu Qiao. Dual relation network for scene text recognition. *IEEE Transactions on Multimedia*, 25:4094–4107, 2023b. doi: 10.1109/TMM.2022.3171108.

Ming Li, Bin Fu, Zhengfu Zhang, and Yu Qiao. Character-aware sampling and rectification for scene text recognition. *IEEE Transactions on Multimedia*, 25:649–661, 2023c. doi: 10.1109/TMM.2021.3129651.

Ming Li, Han Chen, Chenguang Wang, Dang Nguyen, Dianqi Li, and Tianyi Zhou. Ruler: Improving llm controllability by rule-based data recycling, 2024a. URL `https://arxiv.org/abs/2406.15938`.

Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, Jiuxiang Gu, and Tianyi Zhou. Selective reflection-tuning: Student-selected data recycling for LLM instruction-tuning. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 16189–16211, Bangkok, Thailand and virtual meeting, August 2024b. Association for Computational Linguistics. URL `https://aclanthology.org/2024.findings-acl.958`.

Ming Li, Pei Chen, Chenguang Wang, Hongyu Zhao, Yijun Liang, Yupeng Hou, Fuxiao Liu, and Tianyi Zhou. Mosaic it: Enhancing instruction tuning with data mosaics, 2024c. URL `https://arxiv.org/abs/2405.13326`.

Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. Monkey: Image resolution and text label are important things for large multi-modal models. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024d.

Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li, and Wei Peng. A survey on hallucination in large vision-language models. *arXiv preprint arXiv:2402.00253*, 2024a.

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023a.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023b.

Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024b. URL `https://llava-vl.github.io/blog/2024-01-30-llava-next/`.

Michael Xieyang Liu, Frederick Liu, Alexander J. Fiannaca, Terry Koo, Lucas Dixon, Michael Terry, and Carrie J. Cai. "we need structured output": Towards user-centered constraints on large language model output. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, CHI '24. ACM, May 2024c. doi: 10.1145/3613905.3650756. URL `http://dx.doi.org/10.1145/3613905.3650756`.

Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo, Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li. Trustworthy llms: a survey and guideline for evaluating large language models' alignment, 2024d. URL `https://arxiv.org/abs/2308.05374`.

Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. Deepseek-vl: Towards real-world vision-language understanding, 2024.

Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mottaghi, and Aniruddha Kembhavi. Unified-io: A unified model for vision, language, and multi-modal tasks. In *The Eleventh International Conference on Learning Representations*, 2022.

Ruotian Luo and Gregory Shakhnarovich. Comprehension-guided referring expressions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7102–7111, 2017.

Ziqiao Ma, Jiayi Pan, and Joyce Chai. World-to-words: Grounded open vocabulary acquisition through fast mapping in vision-language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 524–544, 2023.

Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 2263–2279, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.177. URL `https://aclanthology.org/2022.findings-acl.177`.

Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. Docvqa: A dataset for vqa on document images. In *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 2199–2208, 2021. doi: 10.1109/WACV48630.2021.00225.

Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and C.V. Jawahar. Infographicvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 1697–1706, January 2022.

Varun K Nagaraja, Vlad I Morariu, and Larry S Davis. Modeling context between objects for referring expression understanding. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pp. 792–807. Springer, 2016.

Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies, 2023. URL `https://arxiv.org/abs/2308.03188`.

Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.

Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, Qixiang Ye, and Furu Wei. Grounding multimodal large language models to the world. In *The Twelfth International Conference on Learning Representations*, 2024.

Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pp. 2641–2649, 2015.

Vipula Rawte, Amit Sheth, and Amitava Das. A survey of hallucination in large foundation models. *arXiv preprint arXiv:2309.05922*, 2023.

Nishant Subramani, Alexandre Matton, Malcolm Greaves, and Adrian Lam. A survey of deep learning approaches for ocr and document understanding, 2021. URL https://arxiv.org/abs/2011.13534.

Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, Zhengliang Liu, Yixin Liu, Yijue Wang, Zhikun Zhang, Bertie Vidgen, Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, Chunyuan Li, Eric Xing, Furong Huang, Hao Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis, Marinka Zitnik, Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei Han, Jieyu Zhao, Jiliang Tang, Jindong Wang, Joaquin Vanschoren, John Mitchell, Kai Shu, Kaidi Xu, Kai-Wei Chang, Lifang He, Lifu Huang, Michael Backes, Neil Zhenqiang Gong, Philip S. Yu, Pin-Yu Chen, Quanquan Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Jana, Tianlong Chen, Tianming Liu, Tianyi Zhou, William Wang, Xiang Li, Xiangliang Zhang, Xiao Wang, Xing Xie, Xun Chen, Xuyu Wang, Yan Liu, Yanfang Ye, Yinzhi Cao, Yong Chen, and Yue Zhao. Trustllm: Trustworthiness in large language models, 2024a. URL https://arxiv.org/abs/2401.05561.

Yu Sun, Dongzhan Zhou, Chen Lin, Conghui He, Wanli Ouyang, and Han-Sen Zhong. Locr: Location-guided transformer for optical character recognition, 2024b. URL https://arxiv.org/abs/2403.02127.

MS Phi-3 Team. Phi-3 technical report: A highly capable language model locally on your phone, 2024. URL https://arxiv.org/abs/2404.14219.

Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature medicine*, 29(8):1930–1940, 2023.

Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pp. 23318–23340. PMLR, 2022.

Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvlm: Visual expert for pretrained language models, 2023a.

Yonghui Wang, Wengang Zhou, Hao Feng, Keyi Zhou, and Houqiang Li. Towards improving document understanding: An exploration on text-grounding via mllms, 2023b. URL https://arxiv.org/abs/2311.13194.

Haoran Wei, Chenglong Liu, Jinyue Chen, Jia Wang, Lingyu Kong, Yanming Xu, Zheng Ge, Liang Zhao, Jianjian Sun, Yuang Peng, Chunrui Han, and Xiangyu Zhang. General ocr theory: Towards ocr-2.0 via a unified end-to-end model, 2024. URL https://arxiv.org/abs/2409.01704.

Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao, and Tianyi Zhou. A survey on knowledge distillation of large language models, 2024. URL https://arxiv.org/abs/2402.13116.

Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Faisal Ahmed, Zicheng Liu, Yumao Lu, and Lijuan Wang. Unitab: Unifying text and box outputs for grounded vision-language modeling. 2022.

Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, Qianyu Chen, Huarong Zhou, Zhensheng Zou, Haoye Zhang, Shengding Hu, Zhi Zheng, Jie Zhou, Jie Cai, Xu Han, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint 2408.01800*, 2024.

Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity. *arXiv preprint arXiv:2310.07704*, 2023.

Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity. In *The Twelfth International Conference on Learning Representations*, 2024. URL `https://openreview.net/forum?id=2msbbX3ydD`.

Licheng Yu, Hao Tan, Mohit Bansal, and Tamara L Berg. A joint speaker-listener-reinforcer model for referring expressions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7282–7290, 2017.

Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024a.

Ruiyi Zhang, Yanzhe Zhang, Jian Chen, Yufan Zhou, Jiuxiang Gu, Changyou Chen, and Tong Sun. Trins: Towards multimodal language models that can read. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 22584–22594, June 2024b.

Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Yu Liu, Kai Chen, and Ping Luo. Gpt4roi: Instruction tuning large language model on region-of-interest. *arXiv preprint arXiv:2307.03601*, 2023a.

Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. LLaVAR: Enhanced visual instruction tuning for text-rich image understanding. *arXiv preprint arXiv:2306.17107*, 2023b.

Yichi Zhang, Ziqiao Ma, Xiaofeng Gao, Suhaila Shakiah, Qiaozi Gao, and Joyce Chai. Groundhog: Grounding large language models to holistic segmentation, 2024c. URL `https://arxiv.org/abs/2402.16846`.

Yang Zhao, Zhijie Lin, Daquan Zhou, Zilong Huang, Jiashi Feng, and Bingyi Kang. Bubogpt: Enabling visual grounding in multi-modal llms. *arXiv preprint arXiv:2307.08581*, 2023.

# A  DETAILED EVALUATION SETTINGS AND METRICS

To simulate various real-world scenarios and test models' capability for grounding at different levels, three different evaluation settings are provided compatible with this benchmark.
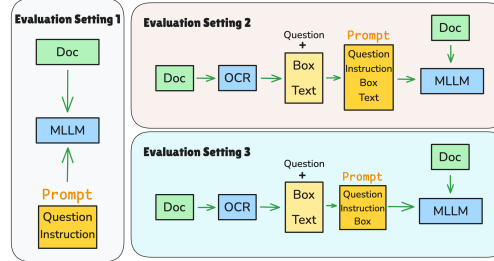


Figure 3: **Illustrations on Evaluation Settings.** In evaluation setting 1, no OCR model is used, representing the hardest scenario. While in settings 2 & 3, an additional OCR model is utilized to facilitate LLM on grounding information generation. The "Instruction" in the prompt describes the requirement of generating grounded bounding boxes and defines the desired format.

## A.1  EVALUATION SETTING 1 (OCR-FREE GROUNDING)

In this setting, as shown in Figure 3 (a), only the document image and corresponding question are provided for the MLLMs, together with the specific instruction that describes the task requirement and defines the desired format. It represents the hardest level of MLLMs' grounding capabilities as it requires MLLMs to answer the question and simultaneously generate the corresponding grounded bounding boxes from scratch. In addition to the measurement of MLLMs' grounding capabilities, it also measures their instruction-following abilities. As a complex customized task, the instructions for generating grounded bounding boxes have never been seen by MLLMs during training, and it requires deep reasoning capability to correctly follow.

**Pixel-level IoU** (Intersection over Union) is utilized as the evaluation metric. It calculates the overlap between pixels in the predicted bounding boxes and ground truth bounding boxes, normalized by the total number of unique pixels. It is commonly used in image segmentation tasks to assess the accuracy of pixel-wise predictions:

$$\text{IoU}_{\text{pixel}} = \frac{1}{N} \sum_{i=1}^{N} \frac{|\text{Pred}_i \cap \text{GroundTruth}_i|}{|\text{Pred}_i \cup \text{GroundTruth}_i|} \tag{1}$$

where $N$ represents the total number of testing samples, $\text{Pred}_i$ represents the pixels within the predicted grounded bounding boxes of $i$th testing sample and $\text{GroundTruth}_i$ represents the pixels within the groundtruth bounding boxes.

## A.2  EVALUATION SETTING 2 (OCR-BASED GROUNDING)

Although the previous setting aligns with our aim the best, it is too hard for all the existing models, from proprietary models to open-source models. Thus further evaluation settings are proposed. In this setting, as shown in Figure 3 (b), an additional OCR model is utilized to facilitate the generation of grounded bounding boxes. Specifically, all the bounding box coordinates and corresponding text content obtained from the OCR model will be wrapped into the prompt for MLLMs as additional information. Under this circumstance, the bounding box generation task is converted into an easier bounding box selection/retrieval task, in which MLLMs do not need to generate coordinates from scratch but only need to select the proper ones given in the prompt. It is worth noting that this evaluation setting still measures MLLMs' grounding capability as they have to align the given OCR information to the image information. In this setting, instance-level IoU score, precision, recall, and F1 score are utilized as the evaluation metrics.

**Instance-level IoU** measures the overlap between the retrieved grounded bounding boxes and the ground truth bounding boxes, normalized by the total number of unique elements. It evaluates how well the retrieved bounding boxes match the ground truth ones:

$$\text{IoU}_{\text{inst}} = \frac{1}{N} \sum_{i=1}^{N} \frac{|\text{Pred}_i \cap \text{GroundTruth}_i|}{|\text{Pred}_i \cup \text{GroundTruth}_i|} \tag{2}$$

where $N$ represents the number of testing samples, $\text{Pred}_i$ represents the bounding boxes selected from the prompt of $i$th testing sample and $\text{GroundTruth}_i$ represents the groundtruth bounding boxes.

**Precision** measures the proportion of correctly retrieved elements out of all elements retrieved by the model. It reflects the accuracy of the model's positive predictions:

$$\text{Precision} = \frac{1}{N} \sum_{i=1}^{N} \frac{|\text{Pred}_i \cap \text{GroundTruth}_i|}{|\text{Pred}_i|} \tag{3}$$

**Recall** measures the proportion of correctly retrieved elements out of all actual ground truth elements. It indicates the model's ability to capture all relevant elements in the ground truth.

$$\text{Recall} = \frac{1}{N} \sum_{i=1}^{N} \frac{|\text{Pred}_i \cap \text{GroundTruth}_i|}{|\text{GroundTruth}_i|} \tag{4}$$

**F1 Score** is the harmonic mean of precision and recall, providing a balanced measure that accounts for both false positives and false negatives.

$$\text{F1} = \frac{1}{N} \sum_{i=1}^{N} \frac{2 \times \text{Precision}_i \times \text{Recall}_i}{\text{Precision}_i + \text{Recall}_i} \tag{5}$$

where $\text{Precision}_i$ and $\text{Recall}_i$ represent the precision and recall of $i$th testing sample.

## A.3 EVALUATION SETTING 3

Although the above two settings represent the most common scenario of grounded bounding box generation for document-level VQA tasks, however, they both have their own flaws: Setting 1 is too difficult for existing MLLMs if not specifically trained on in-domain data, even GPT4o can only have a performance under 10%, thus not suitable for the evaluation of most existing MLLMs. Setting 2, on the contrary, provides a potential shortcut that models might directly select the bounding boxes according to the text content.

Thus, Setting 3, as shown in Figure 3 (c), is proposed as the combination of previous settings, where the OCR-generated bounding boxes will be provided without the specific text content. Under this circumstance, the testing model still needs to have a spatial understanding to select the correct grounded bounding boxes with no need to generate coordinates from scratch. Though this setting is not aligned with the most common practical scenarios, it serves as a reasonable evaluation setting. The evaluation metrics used for this setting are the same as in the setting 2.

## A.4 EVALUATION PROMPTS

In order to alleviate the influences of prompts when testing on our benchmark, 3 different prompts with diverse formatting requirements (Liu et al., 2024c; Li et al., 2024c;a) are proposed for each evaluation setting, resulting in a total of 9 different evaluation prompts. For each evaluation setting, we present the best results across the 3 different prompts as shown in Figure 4. The main difference between different prompts in the same evaluation setting is the required format for the generation or selection of the grounded bounding boxes. Specifically, prompts with "#1" utilize the CSS format, which is widely used for pretraining; prompts with "#2" utilize the most naive format containing the necessary information in a list; prompts with "#3" are similar to "#2" but utilize the relative coordinates.

Prompt # 1 for Evaluation Setting 1

{question}
Please first answer the question according to the given image and then generate the grounded text bounding boxes that support your answer. The width of the image is {width}px and the height is {height}px. Please generate the bounding boxes in the CSS format:
textElement {{
    left: {left}px;
    top: {top}px;
    width: {width}px;
    height: {height}px;
    content: "{text}";
}}

Prompt # 2 for Evaluation Setting 1

{question}
Please first answer the question according to the given image and then generate the grounded text bounding boxes that support your answer. The width of the image is {width}px and the height is {height}px. Please generate the bounding boxes in the below format, the coordinates should be intergers:
[{left}, {top}, {width}, {height}, "{text}"]

Prompt # 3 for Evaluation Setting 1

{question}
Please first answer the question according to the given image and then generate the grounded text bounding boxes that support your answer. Please generate the bounding boxes in the below format using relative coordinates:
[{left}, {top}, {width}, {height}, "{text}"]

Prompt # 1 for Evaluation Setting 2

{question}
Please first answer the question according to the given image and then select the grounded bounding boxes that support your answer. The width of the image is {width}px and the height is {height}px. All the bounding boxes are provided below in the CSS format:
textElement-1 {{
    left: {left}px;
    top: {top}px;
    width: {width}px;
    height: {height}px;
    content: "{text}";
}},
textElement-2 {{
    left: {left}px;
    top: {top}px;
    width: {width}px;
    height: {height}px;
    content: "{text}";
}},
...

Prompt # 2 for Evaluation Setting 2

{question}
Please first answer the question according to the given image and then select the grounded bounding boxes that support your answer. The width of the image is {width}px and the height is {height}px. All the bounding boxes are provided below in the below format::
[1, [{left}, {top}, {width}, {height}, "{text}"]],
[2, [{left}, {top}, {width}, {height}, "{text}"]],
...

Prompt # 3 for Evaluation Setting 2

{question}
Please first answer the question according to the given image and then select the grounded bounding boxes that support your answer. The width of the image is {width}px and the height is {height}px. All the bounding boxes are provided below in the below format::
[1, [{left}, {top}, {width}, {height}, "{text}"]],
[2, [{left}, {top}, {width}, {height}, "{text}"]],
...

Prompt # 1 for Evaluation Setting 3

{question}
Please first answer the question according to the given image and then select the grounded bounding boxes that support your answer. The width of the image is {width}px and the height is {height}px. All the bounding boxes are provided below in the CSS format:
textElement-1 {{
    left: {left}px;
    top: {top}px;
    width: {width}px;
    height: {height}px;
}},
textElement-2 {{
    left: {left}px;
    top: {top}px;
    width: {width}px;
    height: {height}px;
}},
...

Prompt # 2 for Evaluation Setting 3

{question}
Please first answer the question according to the given image and then select the grounded bounding boxes that support your answer. The width of the image is {width}px and the height is {height}px. All the bounding boxes are provided below in the below format::
[1, [{left}, {top}, {width}, {height}]],
[2, [{left}, {top}, {width}, {height}]],
...

Prompt # 3 for Evaluation Setting 3

{question}
Please first answer the question according to the given image and then select the grounded bounding boxes that support your answer. The width of the image is {width}px and the height is {height}px. All the bounding boxes are provided below in the below format::
[1, [{left}, {top}, {width}, {height}]],
[2, [{left}, {top}, {width}, {height}]],
...

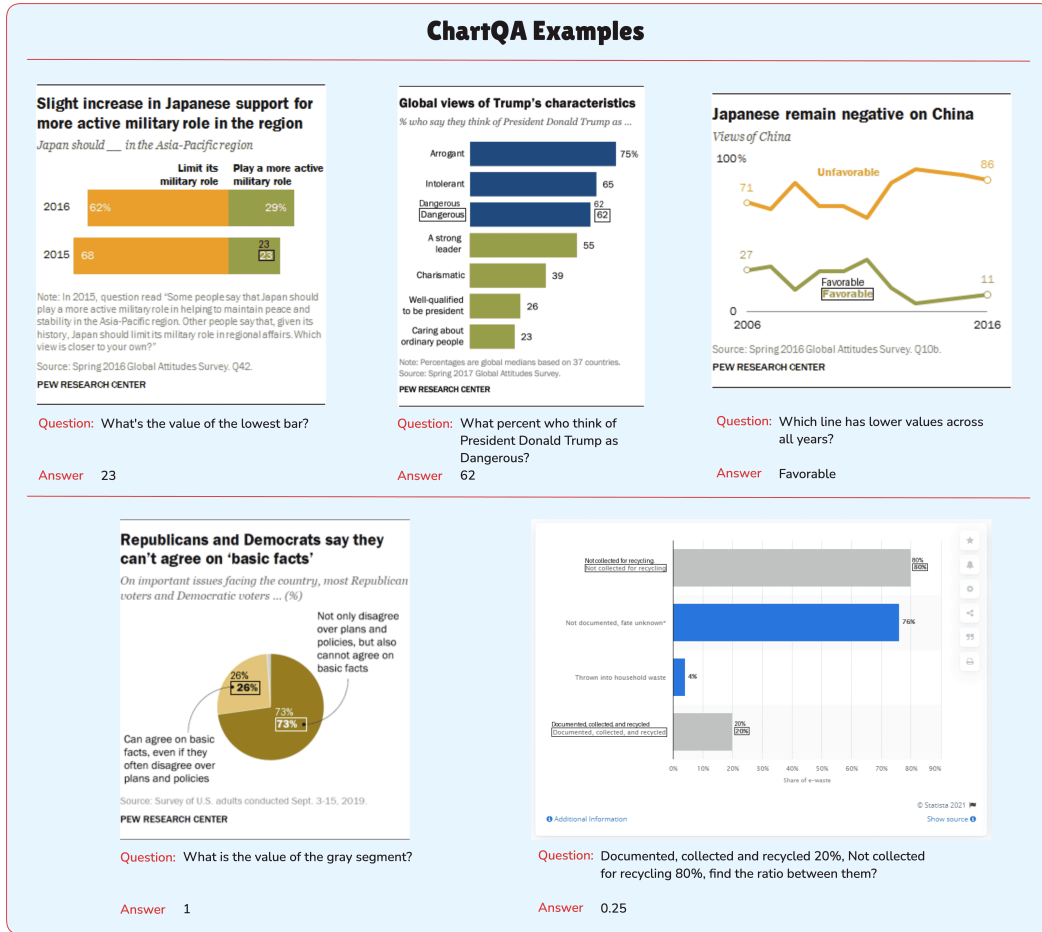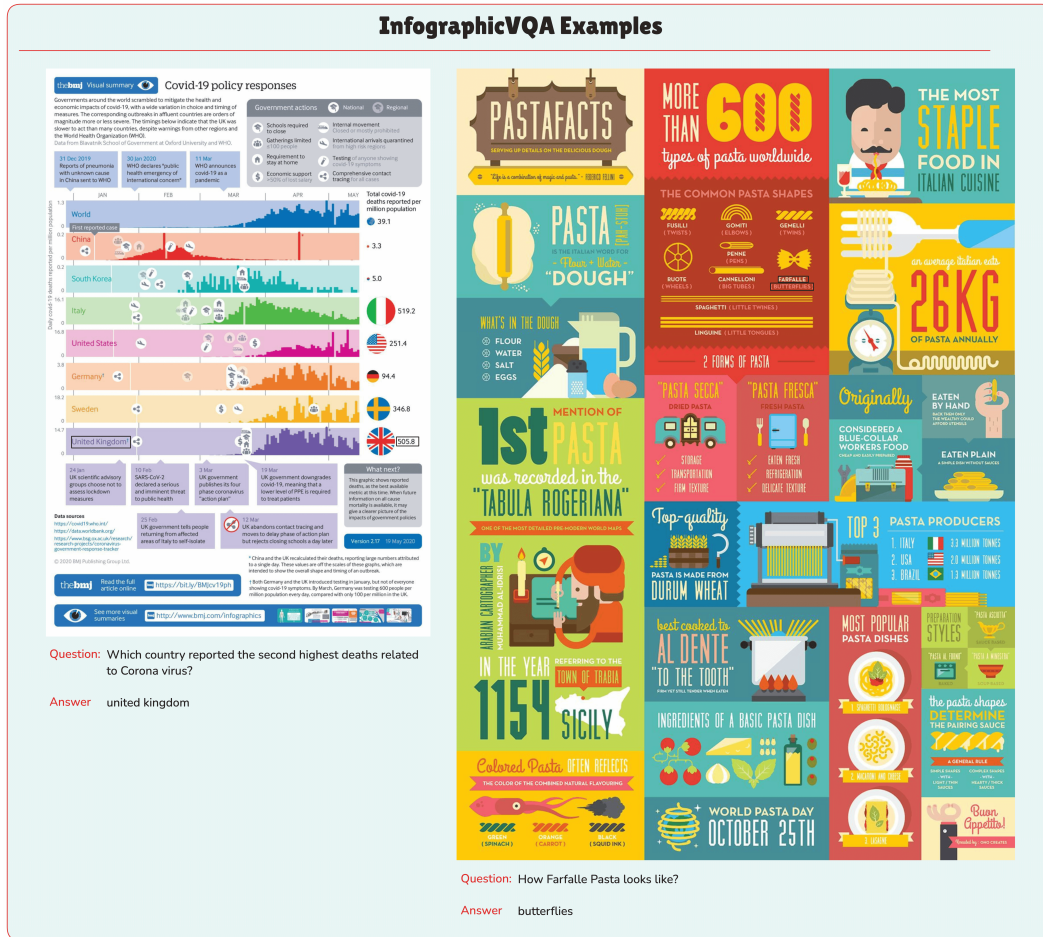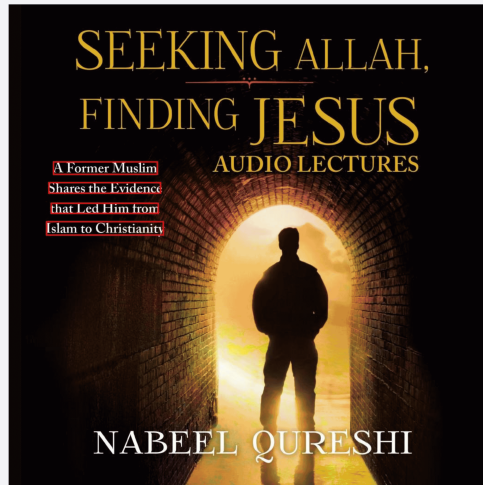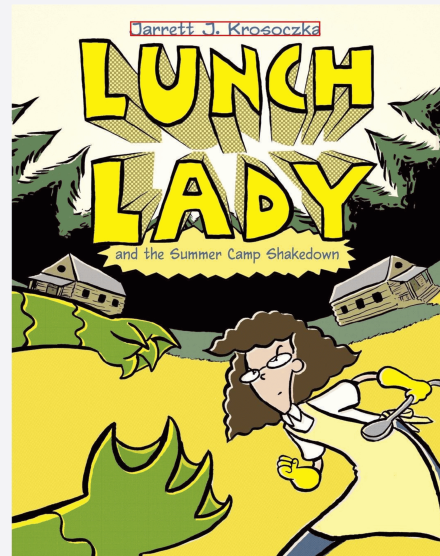Figure 4: The evaluation prompts for different evaluation settings.

Figure 5: Benchmark data examples from ChatQA. The grounded bounding boxes have already been visualized in the original image for better illustration.

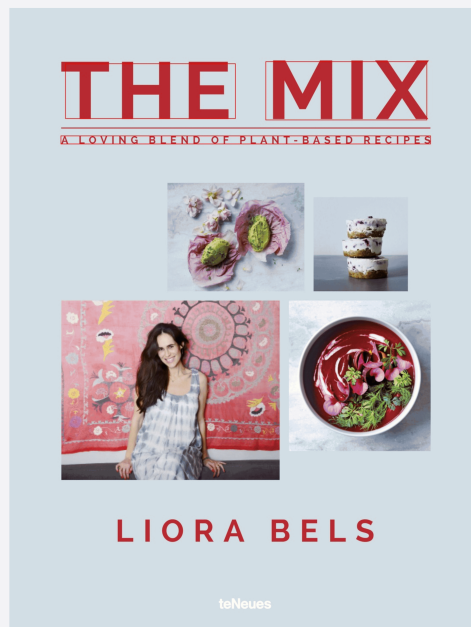## B  BENCHMARK DATA EXAMPLES

The benchmark data examples are visualized in Figure 5 for ChatQA, Figure 6 for DocVQA, Figure 7 for InfographicsVQA, Figure 8 for TRINS. The questions and answers are given in text form, the supported grounded bounding boxes are directly visualized in the images. The samples from different sources show the diversity of our benchmark data.

Figure 6: Benchmark data examples from DocVQA. The grounded bounding boxes have already been visualized in the original image for better illustration.

Figure 7: Benchmark data examples from InfographicsVQA. The grounded bounding boxes have already been visualized in the original image for better illustration.

Figure 8: Benchmark data examples from TRINS. The grounded bounding boxes have already been visualized in the original image for better illustration.