

🐶 DOGE: Towards Versatile Visual Document Grounding and Referring

Yinan Zhou^{1,2,3*}, Yuxin Chen^{2*,†}, Haokun Lin^{2,3,4}, Shuyu Yang¹,
Li Zhu¹, Zhongang Qi^{2‡}, Chen Ma^{3‡}, and Ying Shan²

¹Xi'an Jiaotong University ²ARC Lab, Tencent PCG ³City University of Hongkong ⁴Institute of Automation, CAS

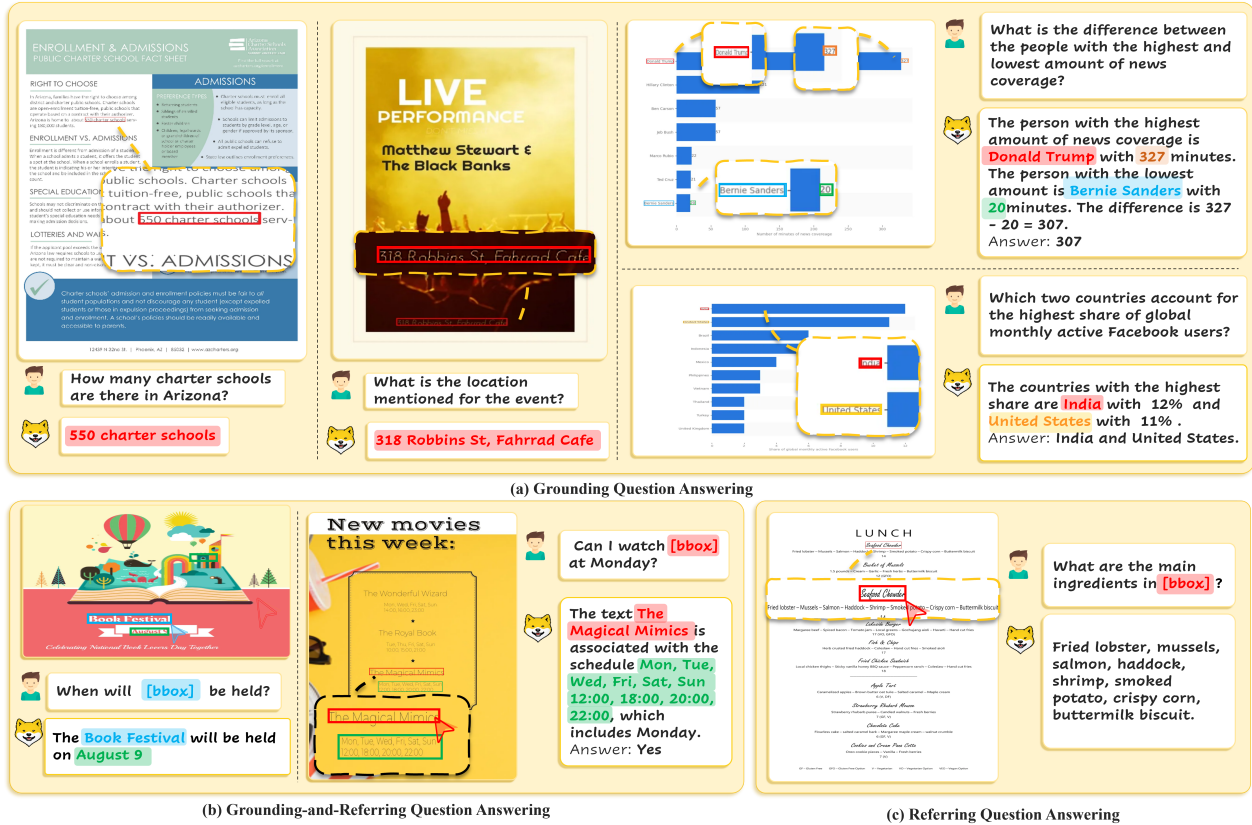


Figure 1. We propose DOGE, a multi-modal large language model that enables users to engage in versatile grounded document interactions. The figure illustrates the reasoning results of DOGE for grounding, grounding-and-referring, and referring tasks in DOGE-Bench.

Abstract

In recent years, Multimodal Large Language Models (MLLMs) have increasingly emphasized grounding and referring capabilities to achieve detailed understanding and flexible user interaction. However, in the realm of visual document understanding, these capabilities lag behind due to the scarcity of fine-grained datasets and comprehensive benchmarks. To fill this gap, we propose the **DO**ocument **G**rounding and **rE**ferring data engine (**DOGE-Engine**), which produces two types of high-quality fine-grained document data: multi-granular parsing data for enhancing fun-

damental text localization and recognition capabilities; and instruction-tuning data to activate MLLM’s grounding and referring capabilities during dialogue and reasoning. Additionally, using our engine, we construct **DOGE-Bench**, which encompasses 7 grounding and referring tasks across 3 document types (chart, poster, PDF document), providing comprehensive evaluations for fine-grained document understanding. Furthermore, leveraging the data generated by our engine, we develop a strong baseline model, **DOGE**. This pioneering MLLM is capable of accurately referring and grounding texts at multiple granularities within document images. Our code, data, and model will be open-sourced for community development.

*Equal contribution.
†Project lead.
‡Corresponding authors.

1. Introduction

In recent years, Multimodal Large Language Models (MLLMs) [6, 14, 21, 23, 25, 26, 32, 33, 62–64] have achieved significant advancements in general visual understanding and reasoning by integrating pre-trained vision encoders with large language models. Leveraging the diverse fine-grained annotations of existing image datasets, some researchers further equip MLLMs with grounding and referring capabilities. These capabilities enhance the detailed visual understanding, increase the credibility of responses, and facilitate more efficient human-AI interaction.

In the field of visual document understanding, the dense textual content and complex layout significantly complicate fine-grained understanding.

To build a user-friendly and trustworthy document AI assistant, it is crucial to enable users to refer to specific regions of a document for more precise comprehension and provide accurate grounding of key details for more expressive and effective interaction.

However, due to the lack of high-quality, fine-grained document datasets, the full potential of MLLMs in document grounding and referring remains largely unexplored. Existing efforts leverage the multi-granularity parsing annotations to enhance document detailed perception [12], or introduce region-level instruction-tuning tasks to achieve basic referring capabilities [28]. However, these works have two significant shortcomings:

- **Suboptimal annotation quality.** To collect parsing annotations, the Optical Character Recognition (OCR) tools are typically employed to extract text and bounding boxes. Subsequently, random document regions and overlapping text boxes are selected to create multi-granular annotations. However, OCR tools have issues with inaccurate bounding box predictions, and the random selection can lead to truncated characters and imprecise bounding boxes. Furthermore, the text annotated using this approach often suffers from incomplete semantics, limiting its potential for application in developing instruction-tuning data.
- **Lack of diversity in task formats.** Current instruction-tuning datasets primarily cater to fundamental referencing tasks, such as region-level OCR and summarization. However, these datasets fall short in supporting grounding tasks and fail to seamlessly incorporate grounding and referencing capabilities into the dialogue and reasoning processes of MLLMs. The limited range of tasks hinders the model’s ability to perceive details flexibly and impacts the user interaction experience.

In order to advance the document referring and grounding capability, we propose the **DO**ocument **G**rounding and **rE**ffering data engine (**DOGE-Engine**) for high-quality fine-grained document data construction. With DOGE-

Engine, we create 1) 1.4M multi-granular document parsing data, which includes text box annotations at the word, phrase, line, paragraph and full-page level across poster, chart and PDF document. This dataset is utilized to enhance basic text localization and recognition capabilities and serves as the foundation for creating instruction-tuning data; 2) a diverse set of 700K instruction-tuning data. This includes both text-in location-out (grounding) and location-in text-out (referring) data, as well as data that combines location and text in both input and output. These instruction-tuning data are constructed based on our multi-granular document parsing data via the assistance of GPT-4o [14], possessing high linguistic quality and accurate ground-and-refer annotations.

Furthermore, in visual document understanding, evaluation tasks related to grounding and referring are relatively scarce, making it challenging to assess the model’s corresponding capabilities. To this end, we propose **DOGE-Bench**, which contains 4K test samples and encompasses 7 grounding and referring tasks across 3 document types (chart, poster, PDF document). Finally, based on data generated by our engine, we develop a strong baseline model, **DOGE**, capable of understanding spatial referring and accurately grounding text within document images. We report the performance of our model on DOGE-Bench, providing a performance reference for future research.

In summary, our contributions are threefold. (1) We introduce DOGE-Engine, a data construction pipeline that generates large-scale, high-quality multi-granular document parsing data and diverse ground-and-refer instruction tuning data. (2) We develop DOGE-Bench, the first comprehensive benchmark designed to evaluate MLLM’s grounding and referring capabilities in document understanding. (3) We present DOGE, a pioneering MLLM that is capable of understanding referred text and performing text grounding during conversation.

2. Related Work

2.1. MLLMs for Visual Document Understanding

Visual Document Understanding focuses on comprehending images with extensive text content. Recently, several Multimodal Large Language Models [6, 9, 56, 61] have been introduced to perform visual document understanding without relying on OCR tools. UReader [61] unifies a wide range of document understanding tasks with instruction-tuning format. A shape-adaptive cropping module is further designed to encode rich textual content in high-resolution image, where the raw image is cropped into multiple tiles, with each tile being individually encoded to represent features of the raw image. TextMonkey [34] employs shifted window to build connections between different image tiles, alleviating the issue of incoherence semantic caused by im-

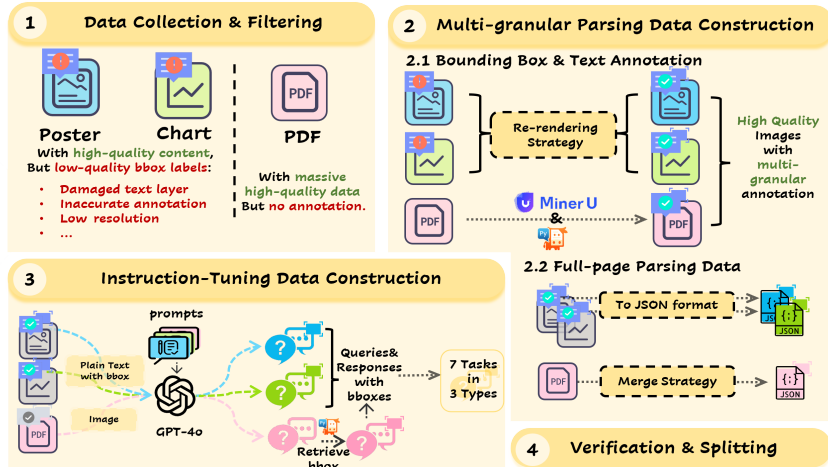
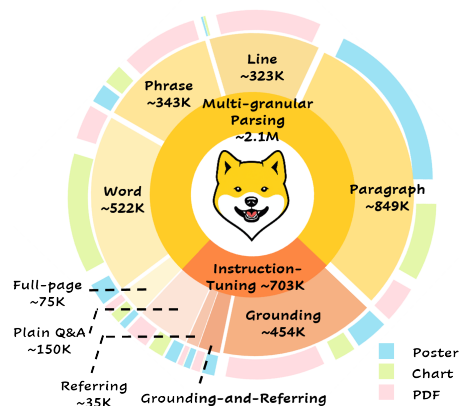


Figure 2. Left: Hierarchical task taxonomy and sample distribution analysis. Right: The pipeline of DOGE-Engine, which outlines the construction process for multi-granularity parsing data and ground-and-refer instruction-tuning data.

age cropping. InternLM-XComposer-4KHD [9] and InternVL 1.5 [7] further increase the tile number, significantly improving the performance on visual document understanding tasks. These works achieve promising performance but lack document grounding and referring capabilities, which hinders the grounded document understanding and flexible human-AI interaction.

2.2. MLLMs for Grounding and Referring

In pursuit of fine-grained image understanding and convenient interaction, recent studies integrate grounding and referring abilities into MLLMs [4, 44, 65, 68]. Kosmos2 [44], Shikra [4] and Ferret [65] utilize bounding boxes or visual prompts to pinpoint specific regions of an image and generate responses with key objects being grounded, facilitating flexible content referring and interaction. Additionally, LLaVA-Grounding [68] and GLaMM [45] further employ finer-grained multi-granularity masks for pixel-level grounding across various semantic levels. These works achieve promising grounding and referring results on real-world images but cannot be adapted to visual document understanding due to the image domain gap. Recently, there are several attempts to develop visual document grounding and referring. mPLUG-1.5 [12] and Kosmos-2.5 [35] enhance localization and fine-grained perception with bounding box annotations from OCR tools and PDF parsing tools, respectively. Fox [29] utilizes various visual prompts to refer to document regions and enables the MLLMs to extract or translate the region-level content. However, mPLUG-1.5 and Kosmos-2.5 only support basic text localization and recognition tasks. Fox further supports the referring translation task.

These methods, though effective, fall short of integrating both grounding and referring into the broader reasoning and dialogue processes. As a result, they leave the full po-

tential of MLLMs for visual document grounding and referring largely untapped, particularly in reasoning tasks where grounding plays a crucial role.

3. DOGE-Engine

Well-annotated and diverse grounded data are crucial for improving the grounding and referring capabilities of MLLMs [27, 44, 52, 67]. Currently, there is a shortage of comprehensive and accurately labeled document grounded data. Manually annotating raw document images is both time-consuming and labor-intensive, as it requires not only marking the bounding boxes but also accurately annotating all the text within those boxes. To tackle this challenge, we collect a substantial volume of documents and develop the **DOGE-Engine** to construct fine-grained document grounded datasets. Our document data sources primarily encompass three document data types: posters, charts, and PDF documents. As shown in the right of Fig. 4, we start by filtering the raw data to remove low-quality samples or those with missing or broken information. Then, we introduce the generation of two types of high-quality document data: The annotation process of multi-granular parsing data is detailed in Sec. 3.1, and the construction pipeline of instruction-tuning data is elaborated in Sec. 3.2. An overview of statistical information is provided in Sec. 3.3. Annotation results are presented in Appendix A.

3.1. Multi-granular Parsing Data Construction

3.1.1. Automatic Bounding box Annotation

Poster. We collect poster data from the Crello dataset [59], which consists of design templates from the design service website. Each poster contains the document meta-annotation, which includes rendered text blocks and their corresponding bounding boxes. However, we observe that

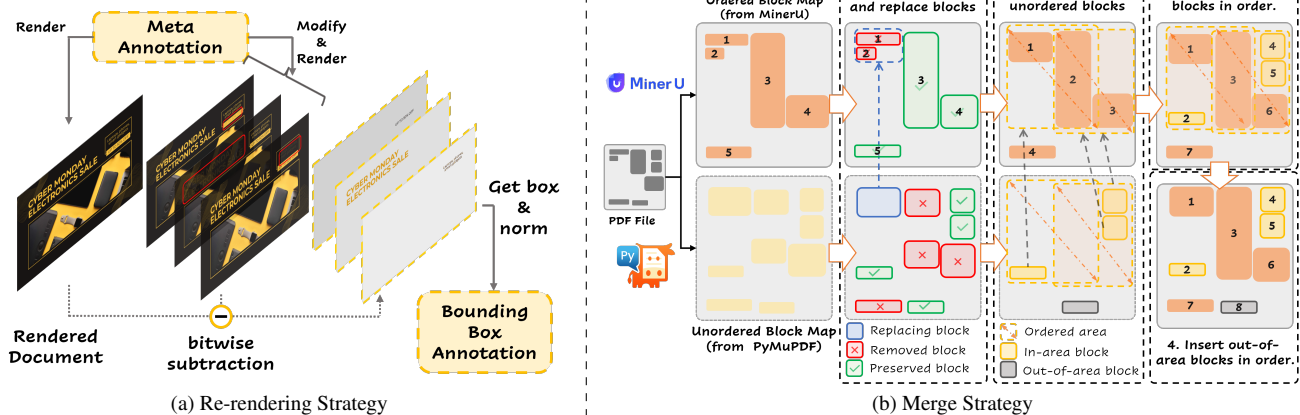


Figure 3. (a) For the poster and chart data, we propose the Re-rendering Strategy to automatically obtain precise bounding boxes. (b) For the PDF document data, we propose the Merge Strategy combining the annotations from MinerU and PyMuPDF to obtain comprehensive and layout-aware full-page parsing annotations

some provided rendered text blocks are significantly degraded, and the bounding boxes are inaccurate. To address this issue, we re-render the text content and design a **Re-rendering Strategy** to obtain precise bounding boxes. Specifically, as illustrated in Fig. 3a, the strategy consists of three steps: Firstly, we render the poster with the meta-annotation; Then, we modify the color or opacity attribute of one text block and perform a re-rendering process; Finally, since the rendering results of first two steps are identical except for the modified attribute, the image of the target text block can be obtained by applying pixel-wise subtraction. By repeating the above steps for every text block, we can obtain accurate bounding boxes of all text blocks, which are then normalized to create the final annotation. Note that we preserve the origin font format during re-rendering to keep the font diversity of poster data.

Chart. The chart data is collected from the ChartQA dataset [38], including various bar charts, line graphs, and pie charts, along with their corresponding JSON/CSV information. We extract the necessary information for rendering images and utilize Matplotlib to create chart images that closely resemble the original charts. To increase the position diversity of text blocks, we apply random padding around the edges of the chart images. Additionally, to prevent the model from over-relying on text and neglecting its ability to interpret visual elements, e.g., bars and lines, we propose to remove the text data from one-third of the chart data, and randomly mask half of the text in another third. Finally, we reuse the Re-rendering Strategy to obtain bounding box annotations for all text blocks, including the values, titles, legends, and axis labels.

PDF Document. CC-MAIN-2021-31-PDF-UNTRUNCATED consists of a large collection of various text documents in PDF format. We select a high-quality subset of moderately sized files. Leveraging the PDF parsing tool PyMuPDF and a document content extraction framework

MinerU[55], we directly extract parsing data of different granularities (word, phrase, line, paragraph).

3.1.2. Full-page Parsing Data Construction

In addition to the detailed annotations across word, phrase, line, and paragraph-level, we also construct full-page parsing annotations to enhance the comprehensive perception of the document content.

Poster and Chart. The poster data features less text, making it straightforward to grasp the dependencies between text blocks. Therefore, we concatenate the paragraph-level text and their corresponding bounding boxes in a left-to-right, top-to-bottom scanning order to create the full-page parsing annotation in JSON format. Due to the highly structured nature of chart data, its full-page parsing annotation is organized as a JSON dictionary containing key-value pairs of chart, axis labels, legends, and their corresponding bounding boxes.

PDF Document. PDF document contains extensive text and complex layouts, necessitating parsing in a logical reading order to comprehend the information accurately. MinerU contains a layout detection model and an OCR model, allowing for extracting text blocks with a certain reading order. However, it often fails to capture all content in documents, missing tables, footnotes, and other elements. In contrast, PyMuPDF can thoroughly extract content from documents, but it does not provide an appropriate reading order. We propose a **Merge Strategy** that combines the strengths of MinerU and PyMuPDF to achieve comprehensive, layout-aware annotations. As shown in Fig. 3b, 1) we first compare text blocks from ordered and unordered maps to eliminate duplicate text blocks. For the truncated blocks (red block 1 and red block 2 in the ordered map), we replace them with the corresponding complete block (blue block in the unordered map) to improve the semantic completeness within blocks; 2) In the ordered map, we construct an or-

dered area if two consecutive blocks are placed from the top left to the bottom right. We classify the preserved blocks in the unordered map into two categories: in-area blocks and out-of-area blocks; 3) We insert the in-area blocks into the ordered area and sequentially update the block order in each ordered area using column-major order; 4) For the out-of-area blocks, we insert them into the ordered map. The order of these blocks is then determined based on their positional relationship with their nearest ordered blocks, following column-major order.

3.2. Instruction-Tuning Data Construction

To seamlessly integrate grounding and referring capabilities in the dialogue and reasoning, it is crucial to construct high-quality instruction-tuning data with accurate grounded text in both query and response. Based on the precise bounding boxes annotations from Sec. 3.1, we leverage GPT-4o [14] to generate diverse-formatted instruction-tuning data tailored to various text granularities, encompassing tasks such as question answering, summarization and reasoning. During construction, we also constrain the length and type of the generated response, resulting in different response categories, e.g., short response, open-ended long response. Subsequently, we add a response format prompt to each query to improve the instruction-following quality of generated data. Prompts for GPT-4o are shown in Appendix B

Poster and Chart. As introduced in Sec. 3.1.2, the full-page parsing annotation for poster and chart data includes all text blocks along with their corresponding bounding boxes, which contain sufficient information to construct ground-and-refer instruction-tuning data. Therefore, we provide GPT-4o with the full-page parsing data, and task GPT-4o to generate queries and responses in a grounded manner, where the generated content must include texts that originate from the given parsing annotations. Additionally, we require GPT-4o to warp the texts originating from parsing annotations with ‘<ocr></ocr>’ and append the corresponding coordinates wrapped in ‘<bbox></bbox>’.

PDF Document. PDF document often contains extensive text content. If we input all the text of various granularities into GPT-4o, it would result in excessively long prompts. This not only incurs significant API calling costs but also increases the difficulty for GPT-4o in comprehending and completing the tasks, leading to low generation quality. To address these issues, we design a **Post-annotating Strategy** to generate high-quality data with minimal overhead. Specifically, 1) We first use the document image as input for GPT-4o instead of the text of document, which significantly reduces the token count; 2) Then, we task GPT-4o to generate queries and responses based on the content of document images. Additionally, any text within the queries and responses that originates from the document image must be wrapped with “<ocr></ocr>”; 3) We extract the texts

wrapped with “<ocr></ocr>” and utilize the PyMuPDF to retrieve the corresponding bounding boxes; 4) Finally, we warp these bounding boxes with ‘<bbox></bbox>’ and insert them back into the generated content after their corresponding texts. For texts that cannot be located by PyMuPDF, we remove wrapped tokens and convert them to plain text. This approach allows the generation of diverse instruction-tuning data at a low cost while ensuring the quality of grounding and referring annotations.

3.3. Data Verification and Splitting

Although we require GPT-4o to generate grounded responses in the format of “<ocr> text </ocr><bbox> x1, y1 ,x2 ,y2 </bbox>”, GPT-4o sometimes fails to follow our requirement, resulting in wrong coordinates format or missing “<bbox></bbox>”. Therefore, we implement a rule-based filter to remove these defective samples. Additionally, for poster and chart data, we extract the grounded text from the generated content and compare them with the full-page parsing annotations, filtering out samples that contain incorrect grounded text. We also performed a detailed categorization of the data. As shown in Fig. 2, DOGE-Dataset includes 1.4M multi-granular parsing data and 700K diverse-formatted instruction-tuning data across three document types: poster, chart and PDF document. The multi-granular parsing data comprises four fine-grained levels (word, phrase, line, paragraph) and a full-page level. These grounded text at different granularity have precise bounding box annotation. The instruction-tuning data comprises four types: grounding, referring, grounding-and-referring, and plain Q&A. The plain Q&A is derived by removing the grounded content from a portion of the grounding data. We construct the plain Q&A subset to enhance the diversity of the data and maintain the traditional document understanding capabilities. Detailed dataset statistics are shown in Appendix D.

4. DOGE-Bench

We introduce the DOGE-Bench for the evaluation of grounding and referring capabilities of MLLMs on visual document understanding tasks.

Task Definition. As shown in Fig. 4, we systematically construct our benchmark by categorizing our data into distinct classes based on both input and output formats. This classification helps in designing clear evaluation metrics. We divide the input formats into two categories based on the presence of bounding boxes: **Grounded Question (GQ)** with bounding boxes, and **Plain-Text Question (PQ)** without bounding boxes. The output formats are categorized into four classes:

- **Grounded Answer(GA):** The response consists of a brief answer accompanied by its corresponding bounding box.

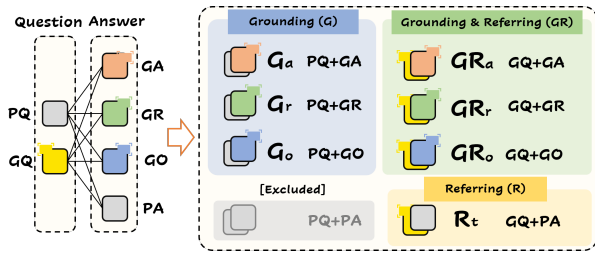


Figure 4. In DOGE-Bench, we categorize the samples into 7 evaluation tasks, further divided into 3 types.

Table 1. Data statistic of DOGE-Bench.

Category	Grounding			Grounding-and-Referring			Referring	Total
	G_a	G_r	G_o	GR_a	GR_r	GR_o	R_t	
#Sample	700	800	636	342	427	775	464	4144

- **Grounded Reasoning (GR):** The response includes the detailed reasoning process and the final answer, while the key text contents in the reasoning process are grounded.
- **Grounded Open-ended Answer (GO):** An open-ended response with one or more key text contents grounded, without providing an answer in a certain format.
- **Plain Text Answer (PA):** This format does not incorporate grounded text content.

By combining two input forms and four output forms, we derive seven document referring and grounding tasks. Among these tasks, three sub-tasks primarily assess grounding capability: grounded answering for plain text questions (G_a), grounded reasoning for plain-text questions (G_r), and grounded open-ended answering for plain-text questions (G_o). The plain-text answering for grounded questions (R_t) task evaluates referring capability. The remaining tasks, grounded answering with plain-text questions (GR_a), grounded reasoning for grounded questions (GR_r), and grounded open-ended answering for grounded questions (GR_o) require the integration of both grounding and referring capabilities for successful completion.

Metrics. Our benchmark evaluation encompasses two aspects: grounding performance and text answer accuracy. Following the previous works in grounded captioning [65], we evaluate the grounding and text answer separately. For grounding performance, we use $F1_{all}$ score, which evaluates grounding results as a multi-label classification problem. The generated grounded text is considered correct if Intersection over Union (IoU) between its bounding box and the GT bounding box is greater than 0.5, meanwhile its text matches with the GT text. For text answer accuracy, we use exact text matching accuracy for short-answer tasks and BLEU score for long-answer tasks.

Data Statistics. Our DOGE-Bench includes 2K grounding samples, 0.5K referring samples, and 15K grounding-and-referring samples. The detailed statistics is shown in Tab. 1.

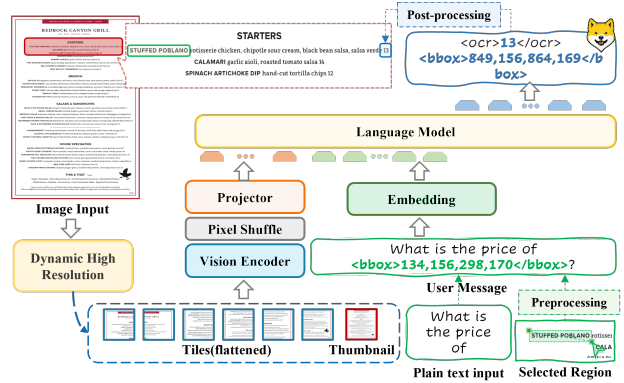


Figure 5. The overall architecture of DOGE.

5. DOGE

Overall Architecture. As illustrated in Fig. 5, our model DOGE employ a general MLLM architecture, including a vision encoder, a projector, and a large language model. In the visual encoder component, to enable the model to handle high resolution, we first search for the best aspect ratio for the input image and dynamically segment images into multiple tiles. These tiles, along with a thumbnail of the input image, are provided as input to the vision encoder. We employ pixel shuffle[7] to improve the computational efficiency of the model when processing high-resolution images. For bounding box representation, we simply discrete the continuous coordinates into discrete values from 0 to 999, avoiding the introduction of extra module or location tokens. During inference, we transfer the user-selected regions to coordinates of bounding boxes and insert them into the query for preprocessing. After obtaining the output, we utilize post-processing to overlay bounding boxes on original document images, thereby facilitating user interaction.

Training Strategy. We adopt a three-stage training strategy, including pre-aligning, pre-training, and fine-tuning. The pre-aligning stage focuses on aligning the feature space of vision and language rapidly. In this stage, we freeze both vision encoder and large language model, and train the projector with a relatively large learning rate. The pre-training stage aims at document parsing capabilities. We unfreeze the vision encoder and the LLM, enabling the model to recognize diverse textual content and acquire text-reading capability. In the fine-tuning stage, We train the entire model using diverse instruction-tuning data, enhancing its instruction-following ability while activating its grounding and referring capabilities during dialogue and reasoning.

Training Dataset. In the pre-aligning stage, we utilize LLaVA-558K [33] to train the projector. For the pre-training dataset, we utilize DocStruct4M [12] along with our 1.4M multi-granular parsing data to enhance basic text reading, text grounding and text referring capabilities of DOGE. For the fine-tuning data, we adopt our ground-and-refer instruction-tuning data and meticulously selected

Table 2. The performance evaluation of DOGE on DOGE-Bench across grounding, grounding-and-referring, and referring conducted for three image categories: Poster, Chart, and PDF Document. Acc is for text accuracy while $F1_{all}$ is for grounding performance.

Document Type	Grounding						Grounding-and-Referring						Referring
	G_a		G_r		G_o		GR_a		GR_r		GR_o		R_t
	Acc	$F1_{all}$	Acc	$F1_{all}$	BLEU4	$F1_{all}$	Acc	$F1_{all}$	Acc	$F1_{all}$	BLEU4	$F1_{all}$	Acc
Poster	82.5	77.5	72	65.7	43	41.6	87.5	85.8	74.5	56.6	45.6	25.8	72.5
Chart	91.5	74	70.5	43.9	32.5	46.1	79.3	60.3	65.9	57.6	39.5	54.4	61
PDF Document	80.3	71	65	45	35.9	29.1	72.6	50	67.3	58.1	37.3	14.7	19.1
Average	84.8	74.2	69.2	51.5	37.1	41.9	79.8	65.4	69.2	57.4	40.8	31.6	50.9

Table 3. Performance comparison on 10 general document benchmarks.

Model	Size	Doc VQA	Info VQA	Deep Form	KLC	WTQ	Tab Fact	Chart QA	Text VQA	Text Caps	Visual MRC
Donut[19]	<1B	67.5	11.6	61.6	30.0	18.8	54.6	41.8	43.5	74.4	93.91
DocOwl[11]	7B	62.2	38.2	42.6	30.3	26.9	60.2	57.4	52.6	111.9	188.8
UReader[61]	7B	65.4	42.2	49.5	32.8	29.4	67.6	59.3	57.6	118.4	221.7
TextMonkey[24]	9B	73.0	28.6	59.7	37.8	31.9	-	66.9	65.9	-	-
Vary[56]	7B	76.3	-	-	-	-	-	66.1	-	-	-
TokenPacker[22]	13B	70.0	-	-	-	-	-	-	-	-	-
DocPeida[10]	7B	47.1	15.2	-	-	-	-	46.9	60.2	-	-
QwenVL[1]	9B	65.1	35.4	-	-	-	-	65.7	63.8	-	-
IXC 2.5[69]	7B	90.9	70.0	71.2	-	<u>53.6</u>	85.2	82.2	78.2	-	307.5
DocOwl-1.5-Chat[12]	8B	82.2	50.7	68.8	<u>38.7</u>	40.6	80.2	70.2	68.6	131.6	246.4
DocOwl-2[13]	8B	80.7	46.4	66.8	37.5	36.5	78.2	70.0	66.7	131.8	217.4
InternVL2[7]	8B	<u>91.6</u>	74.8	-	-	-	-	<u>83.3</u>	<u>77.4</u>	-	-
DOGE	8B	91.7	<u>70.7</u>	<u>70.8</u>	40.4	58.8	<u>84.5</u>	83.6	<u>76.6</u>	145.9	332.5

datasets to enhance model performance across a wide range of document parsing and understanding tasks, resulting in a final fine-tuning dataset of 2M samples. The detailed composition of fine-tuning dataset can be found in supplementary materials. In Appendix E, we provide a detailed description and statistics of our training data composition.

6. Experiment

6.1. Implementation Details

DOGE utilizes the InternViT-300M-448px [6, 7] as the vision encoder and Qwen2-7B-Instruct [60] as the LLM. In the pre-aligning stage, we only train the projector and the learning rate is set to 1e-3. In the pre-training and fine-tuning stage, all model parameters are trainable. The learning rate for the vision encoder is 2e-6, while the learning rate for other components is 1e-5. Each stage is conducted for 1 epoch. In the pre-aligning stage, the batch size is set to 256, and we only use the thumbnail for vision encoder input. In the pre-training stage, the batch size is configured to 512. The number of image tiles and the max length of the input sequence to the large language model are set to 9 and 4096, respectively. In the fine-tuning stage, the batch size is adjusted to 256, and the image tile number and max input length increase to 16 and 6144. During inference, we set the image tile number to 20 for traditional understanding tasks and 16 for grounding and referring tasks.

6.2. Results on DOGE-Bench

Due to the scarcity of models that support document grounding and referring Q&A, we mainly present the performance of our model on DOGE-Bench, providing a performance reference for future research. We evaluate the grounding and referring capabilities of DOGE across three data types and seven tasks.

Posters have a diverse range of font types and colors, making them ideal for assessing a model’s adaptability to different font styles. DOGE shows strong performance in both grounding and referring tasks on poster data, validating its robustness in recognizing fonts with varied styles. Charts contain fine-grained and structural elements, necessitating precise text localization and understanding capabilities. Despite these challenges, DOGE still demonstrates considerable performance. PDF document data typically features high-resolution and complex content, posing challenges on both referring and grounding tasks. Although DOGE performs well in simple question answering (G_a and GR_a) and reasoning (G_r and GR_r), its $F1_{all}$ score in open-ended questions is relatively low. This can be attributed to the difficulty of text localization in dense text images. Some texts could not be accurately located and their bounding boxes are missed during inference, limiting the model’s ability to provide comprehensive grounded outputs. Additional quantitative results are shown in Appendix C.

Table 4. Performance comparison of multi-grained text localization on DocLocal4K dataset.

Model	Text Grounding					Text Recognition				
	word	phrase	line	paragraph	ALL	word	phrase	line	paragraph	ALL
DocOwl-1.5	70.42	76.38	85.88	91.34	80.38	70.10	67.86	73.88	70.70	70.63
DOGE ⁻	74.79	78.41	84.1	95.46	82.64	79.87	75.73	73.88	70.1	75.83
DOGE	81.85	83.58	86.88	95.67	86.64	80.00	78.85	77.60	72.72	77.88

Table 5. The impact of varying the number of tiles during training and inference. #TT denotes the maximum number of tiles during training, and #IT indicates the maximum number during inference.

#TT	#IT	InfoVQA	PDF Document	
		ANLS	Acc_{avg}	$F1_{all_{avg}}$
9	9	68.91	52.88	36.21
16	16	69.75	67.84	42.4
16	20	70.65	50.73	21.0

6.3. Traditional Document Understanding

To assess the overall capability of DOGE, we conduct experiments on 10 traditional document understanding benchmarks, including DocVQA [39], InfographicVQA [40], DeepForm [49] and KLC [48] for document comprehension, WTQ [43] and TabFact [5] for table understanding, ChartQA [37] for chart comprehension, TextVQA [47] and TextCaps [46] for natural image interpretation, VisualMRC [50] for webpage understanding. For metrics, we use ANLS [2] for DocVQA and InfoVQA, F1 score for DeepForm and KLC, text-matching accuracy for WTQ, TabFact, TextVQA and ChartQA, CIDEr [53] for TextCaps and VisualMRC.

As detailed in Tab. 3, we observe that DOGE achieves competitive performance across all tasks. Notably, DOGE outperforms the previous state-of-the-art model IXC2.5 by +5.3% in accuracy on WTQ and +25.0% in CIDEr on VisualMRC. It is important to highlight that our primary focus is on enhancing the grounding and referencing capabilities, and we do not engage in extensive pre-training and fine-tuning on large datasets as InternVL2 [7] and IXC2.5 [69]. Despite the limited data used for training, DOGE ranks second on four datasets while remaining highly competitive, showing that DOGE maintains a strong performance across general document understanding tasks. We believe that incorporating more data for pre-training and fine-tuning could further improve DOGE’s performance on these tasks.

6.4. Ablation Study

Effectiveness of multi-granular parsing data. We evaluate the text recognition and grounding performance on DocLocal4K [12]. As shown in Tab. 4, DOGE surpasses mPLUG-DocOwl-1.5 across all granularities. We also train a model called DOGE⁻, which is pre-trained without using our constructed multi-granular parsing data. Despite the domain gap between our parsing data and DocLocal4K due to their different construction methods, it is evident that incor-

porating our data improves text grounding and recognition accuracy across various text granularities. This validates the effectiveness of our multi-granular parsing data in enhancing basic text grounding and recognition capabilities.

Increasing the number of training image tiles leads to a significant improvement in high-resolution document understanding.

As shown in rows 1 and 2 of Tab. 5, we evaluate model performance with varying numbers of training image tiles on two datasets: InfoVQA and the PDF Document subset of DOGE-Bench, both of which consist of high-resolution images. The results show that increasing the number of training image tiles from 9 to 16 significantly improves answer accuracy on both datasets. This validates the importance of providing sufficiently high-resolution input for effective document understanding in MLLMs.

Increasing the number of image tiles during inference is detrimental to grounding and referring.

As shown in rows 2 and 3 of Tab. 5, we compare model performance with different numbers of image tiles during inference. The results show that increasing the number of tiles improves performance on general document QA tasks, with a 0.9% gain on InfoVQA. However, there is a notable decline in performance on the PDF Document subset of DOGE-Bench, including reduced accuracy for text answers and F1 scores for grounding performance. This suggests that the current bounding box modeling method is sensitive to out-of-domain input resolutions, highlighting the need for a more robust and generalizable bounding box modeling method for varying input resolutions.

7. Conclusion

In this paper, we introduce DOGE-Engine, a data construction pipeline for generating high-quality ground-and-refer data for fine-grained document understanding. Additionally, we construct DOGE-Bench as the first comprehensive benchmark for evaluating the document grounding and referring capabilities of MLLMs. Furthermore, leveraging data generated by our engine, we develop DOGE, a pioneering MLLM that integrates text grounding and referring abilities into the dialogue and reasoning process. Our results show that DOGE can achieve versatile document grounding and referring while achieving promising performance on traditional document understanding tasks. We hope our work will facilitate practical document AI assistants.

References

- [1] 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*, 2023. 7
- [2] Ali Furkan Biten, Rubèn Tito, Andrés Mafla, Lluís Gomez, Marçal Rusiñol, C.V. Jawahar, Ernest Valveny, and Dimosthenis Karatzas. Scene text visual question answering. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4290–4300, 2019. 8
- [3] Jimmy Carter. Textocr-gpt4v. <https://huggingface.co/datasets/jimmycarter/textocr-gpt4v>, 2024. 8
- [4] 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*, 2023. 3
- [5] Wenhui Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyong Zhou, and William Yang Wang. Tabfact : A large-scale dataset for table-based fact verification. In *International Conference on Learning Representations (ICLR)*, Addis Ababa, Ethiopia, 2020. 8
- [6] Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*, 2023. 2, 7, 8
- [7] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. 3, 6, 7, 8
- [8] Zhoujun Cheng, Haoyu Dong, Zhiruo Wang, Ran Jia, Jiaqi Guo, Yan Gao, Shi Han, Jian-Guang Lou, and Dongmei Zhang. Hitab: A hierarchical table dataset for question answering and natural language generation. In *ACL*, 2022. 8
- [9] 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, Wenhui 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*, 2024. 2, 3
- [10] Hao Feng, Qi Liu, Hao Liu, Wengang Zhou, Houqiang Li, and Can Huang. Docpedia: Unleashing the power of large multimodal model in the frequency domain for versatile document understanding, 2023. 7
- [11] Anwen Hu, Yaya Shi, Haiyang Xu, Jiabo Ye, Qinghao Ye, Ming Yan, Chenliang Li, Qi Qian, Ji Zhang, and Fei Huang. mplug-paperowl: Scientific diagram analysis with the multimodal large language model. *arXiv preprint arXiv:2311.18248*, 2023. 7
- [12] Anwen Hu, Haiyang Xu, Jiabo Ye, Ming Yan, Liang Zhang, Bo Zhang, Chen Li, Ji Zhang, Qin Jin, Fei Huang, et al. mplug-docowl 1.5: Unified structure learning for ocr-free document understanding. *arXiv preprint arXiv:2403.12895*, 2024. 2, 3, 6, 7, 8, 9
- [13] Anwen Hu, Haiyang Xu, Liang Zhang, Jiabo Ye, Ming Yan, Ji Zhang, Qin Jin, Fei Huang, and Jingren Zhou. mplug-docowl2: High-resolution compressing for ocr-free multi-page document understanding, 2024. 7
- [14] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 2, 5
- [15] Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. Dvqa: Understanding data visualizations via question answering. In *CVPR*, 2018. 8
- [16] Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Akos Kadar, Adam Trischler, and Yoshua Bengio. Figureqa: An annotated figure dataset for visual reasoning, 2018. 8
- [17] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *ECCV*, 2016. 8
- [18] Aniruddha Kembhavi, Minjoon Seo, Dustin Schwenk, Jonghyun Choi, Ali Farhadi, and Hannaneh Hajishirzi. Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4999–5007, 2017. 8
- [19] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoon Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document understanding transformer. In *Computer Vision – ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVIII*, page 498–517, Berlin, Heidelberg, 2022. Springer-Verlag. 7
- [20] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoon Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document understanding transformer. In *European Conference on Computer Vision (ECCV)*, 2022. 8
- [21] Wen Li, Lixin Duan, Dong Xu, and Ivor Wai-Hung Tsang. Text-based image retrieval using progressive multi-instance learning. In *2011 international conference on computer vision*, pages 2049–2055. IEEE, 2011. 2
- [22] Wentong Li, Yuqian Yuan, Jian Liu, Dongqi Tang, Song Wang, Jie Qin, Jianke Zhu, and Lei Zhang. Tokenpacker: Efficient visual projector for multimodal llm, 2024. 7
- [23] 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*, 2024. 2
- [24] 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*, pages 26763–26773, 2024. 7
- [25] Haokun Lin, Haoli Bai, Zhili Liu, Lu Hou, Muyi Sun, Linqi Song, Ying Wei, and Zhenan Sun. Mope-clip: Structured pruning for efficient vision-language models with module-wise pruning error metric. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27370–27380, 2024. 2
- [26] Haokun Lin, Haobo Xu, Yichen Wu, Jingzhi Cui, Yingtao Zhang, Linzhan Mou, Linqi Song, Zhenan Sun, and Ying Wei. Duquant: Distributing outliers via dual transformation makes stronger quantized llms. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 2
- [27] Weifeng Lin, Xinyu Wei, Ruichuan An, Peng Gao, Bocheng Zou, Yulin Luo, Siyuan Huang, Shanghang Zhang, and Hongsheng Li. Draw-and-understand: Leveraging visual prompts to enable mllms to comprehend what you want, 2024. 3
- [28] Chenglong Liu, Haoran Wei, Jinyue Chen, Lingyu Kong, Zheng Ge, Zining Zhu, Liang Zhao, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Focus anywhere for fine-grained multi-page document understanding. *arXiv preprint arXiv:2405.14295*, 2024. 2
- [29] Chenglong Liu, Haoran Wei, Jinyue Chen, Lingyu Kong, Zheng Ge, Zining Zhu, Liang Zhao, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Focus anywhere for fine-grained multi-page document understanding. *arXiv preprint arXiv:2405.14295*, 2024. 3
- [30] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023. 8
- [31] Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen, Kaiqiang Song, Sangwoo Cho, Yaser Yacoob, and Dong Yu. Mmc: Advancing multimodal chart understanding with large-scale instruction tuning. *arXiv preprint arXiv:2311.10774*, 2023. 8
- [32] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 2
- [33] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 2, 6
- [34] Yulian Liu, Biao Yang, Qiang Liu, Zhang Li, Zhiyin Ma, Shuo Zhang, and Xiang Bai. Textmonkey: An ocr-free large multimodal model for understanding document, 2024. 2
- [35] Tengchao Lv, Yupan Huang, Jingye Chen, Lei Cui, Shuming Ma, Yaoyao Chang, Shaohan Huang, Wenhui Wang, Li Dong, Weiyao Luo, et al. Kosmos-2.5: A multimodal literate model. *arXiv preprint arXiv:2309.11419*, 2023. 3
- [36] U-V Marti and Horst Bunke. The iam-database: an english sentence database for offline handwriting recognition. *International journal on document analysis and recognition*, 5: 39–46, 2002. 8
- [37] 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 *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263–2279, Dublin, Ireland, 2022. Association for Computational Linguistics. 8
- [38] 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 *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263–2279, 2022. 4
- [39] Minesh Mathew, Dimosthenis Karatzas, and C.V. Jawahar. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2200–2209, 2021. 8
- [40] Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and C. V. Jawahar. Infographicvqa. In *2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2582–2591, 2022. 8
- [41] A. Mishra, K. Alahari, and C. V. Jawahar. Scene text recognition using higher order language priors. In *BMVC*, 2012. 8
- [42] Jason Obeid and Enamul Hoque. Chart-to-text: Generating natural language descriptions for charts by adapting the transformer model, 2020. 8
- [43] Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1470–1480, Beijing, China, 2015. Association for Computational Linguistics. 8
- [44] 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, abs/2306*, 2023. 3
- [45] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M. Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S. Khan. Glamm: Pixel grounding large multimodal model. *The IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 3
- [46] Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for image captioning with reading comprehension. In *European Conference on Computer Vision*, 2020. 8
- [47] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8309–8318, 2019. 8
- [48] Tomasz Stanisławek, Filip Graliński, Anna Wróblewska, Dawid Lipiński, Agnieszka Kaliska, Paulina Rosalska, Bartosz Topolski, and Przemysław Biecek. Kleister: Key information extraction datasets involving long documents with complex layouts. In *Document Analysis and Recognition – ICDAR 2021*, pages 564–579, Cham, 2021. Springer International Publishing. 8
- [49] S Svetlichnaya. Deepform: Understand structured documents at scale. 2020. 8

- [50] Ryota Tanaka, Kyosuke Nishida, and Sen Yoshida. Visualmrc: Machine reading comprehension on document images. In *AAAI*, 2021. 8
- [51] Benny J. Tang, Angie Boggust, and Arvind Satyanarayan. Vistext: A benchmark for semantically rich chart captioning, 2023. 8
- [52] Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *arXiv preprint arXiv:2406.16860*, 2024. 3
- [53] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4566–4575, 2015. 8
- [54] Bryan Wang, Gang Li, Xin Zhou, Zhouong Chen, Tovi Grossman, and Yang Li. Screen2words: Automatic mobile ui summarization with multimodal learning, 2021. 8
- [55] Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, Bo Zhang, Liqun Wei, Zhihao Sui, Wei Li, Botian Shi, Yu Qiao, Dahua Lin, and Conghui He. Mineru: An open-source solution for precise document content extraction, 2024. 4
- [56] Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, Jinrong Yang, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Vary: Scaling up the vision vocabulary for large vision-language models. *arXiv preprint arXiv:2312.06109*, 2023. 2, 7
- [57] Chris Wendler. wendlerc/renderedtext, 2023. 8
- [58] Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing. *ArXiv*, abs/2406.08464, 2024. 8
- [59] Kota Yamaguchi. Canvasvae: Learning to generate vector graphic documents. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5481–5489, 2021. 3
- [60] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024. 7
- [61] Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Guohai Xu, Chenliang Li, Junfeng Tian, Qi Qian, Ji Zhang, et al. Ureader: Universal ocr-free visually-situated language understanding with multimodal large language model. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2841–2858, 2023. 2, 7
- [62] Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl3: Towards long image-sequence understanding in multimodal large language models, 2024. 2
- [63] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chaoya Jiang, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. mplug-owl: Modularization empowers large language models with multimodality, 2023.
- [64] Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration, 2023. 2
- [65] 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. 3, 6
- [66] Ye Yuan, Xiao Liu, Wondimu Dikubab, Hui Liu, Zhilong Ji, Zhongqin Wu, and Xiang Bai. Syntax-aware network for handwritten mathematical expression recognition. *arXiv preprint arXiv:2203.01601*, 2022. 8
- [67] Haotian Zhang, Mingfei Gao, Zhe Gan, Philipp Dufter, Nina Wenzel, Forrest Huang, Dhruvi Shah, Xianzhi Du, Bowen Zhang, Yanghao Li, et al. Mm1. 5: Methods, analysis & insights from multimodal llm fine-tuning. *arXiv preprint arXiv:2409.20566*, 2024. 3
- [68] Hao Zhang, Hongyang Li, Feng Li, Tianhe Ren, Xueyan Zou, Shilong Liu, Shijia Huang, Jianfeng Gao, Chunyuan Li, Jainwei Yang, et al. Llava-grounding: Grounded visual chat with large multimodal models. In *European Conference on Computer Vision*, pages 19–35. Springer, 2025. 3
- [69] Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong Duan, Bin Wang, Linke Ouyang, Songyang Zhang, Wenwei Zhang, Yining Li, Yang Gao, Peng Sun, Xinyue Zhang, Wei Li, Jingwen Li, Wenhai Wang, Hang Yan, Conghui He, Xingcheng Zhang, Kai Chen, Jifeng Dai, Yu Qiao, Dahua Lin, and Jiaqi Wang. Internlm-xcomposer-2.5: A versatile large vision language model supporting long-contextual input and output. *arXiv preprint arXiv:2407.03320*, 2024. 7, 8
- [70] Yilun Zhao, Chen Zhao, Linyong Nan, Zhenting Qi, Wenlin Zhang, Xiangru Tang, Boyu Mi, and Dragomir Radev. RobuT: A systematic study of table QA robustness against human-annotated adversarial perturbations. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6064–6081, Toronto, Canada, 2023. Association for Computational Linguistics. 8

A. The Annotation Results

A.1. Poster and Chart Annotations

As shown in Fig. 6, we present a comparison between our annotations and the original annotations.

In the original poster annotations, some text image layers are damaged, and some bounding boxes are inaccurate, which hinders precise text recognition and localization. In the annotations obtained using our Re-rendering Strategy, we reconstruct the text layers and achieve accurate bounding box annotations. Additionally, we include global content information from the original dataset, such as text format, title and keywords to enhance global awareness during instruction-tuning data generation. We use the values of “text_with_box” as the annotations for poster’s full-page parsing data.

In original chart annotations from ChartQA, some bounding boxes are associated with bars/lines rather than text value. This is inconsistent with our text and bounding box correspondence objectives. Additionally, the bounding boxes in the original annotations are not accurate. In our re-rendered annotations, we align the bounding boxes with the text value, and the bounding boxes are accurate. Furthermore, we randomly erase some of the values to ensure that the model can infer the missing values based on other visual information. We use the entire chart JSON dict as the annotations for the full-page parsing of the chart. Due to space constraints, we omit some of the content using ‘...’.

A.2. PDF Document Annotations

As shown in Fig. 7, we present a comparison between the ordered annotations from MinerU and the unordered but comprehensive annotations from PyMuPDF, along with the combined annotations using our Merge Strategy. We utilize green arrows to indicate the ordered annotations and gray arrows to indicate the naive scanning order. By combining the two annotation methods, we achieve full-page parsing annotations that are both comprehensive and as ordered as possible. Moreover, our method can become more effective as the performance of the ordered annotation tools improves. Due to space constraints, we omit the content in the middle of this passage.

B. Prompts and Instructions

B.1. Prompt Details

In Fig. 8, we show three different prompts for poster, chart and PDF document:

For poster, we deploy plain text input as prompt. Besides format information and rule information, we also provide GPT-4o with some of the overall content and style information from the original Crello dataset. This helps in achieving

a better global understanding, thereby improving the generation of instruction-tuning data.

For chart, since the content of the charts only contains some numbers and lacks an introduction to the meaning of the content being statistically represented, we add a question answering data from the original ChartQA to help GPT-4o better understand the meaning conveyed by the chart content. Additionally, when the model generates output, we output the masked values in the grounded format “<ocr>text</ocr><bbox>null</bbox>” as well. After obtaining the output, we perform format filtering to degrade this part of the content into plain text and remove the degraded plain text item in “necessary bbox”.

For PDF document, we directly send the images and simple output format rules to GPT-4o, obtaining the output with the original text wrapped in “<ocr></ocr>”. Then, we use PyMuPDF to query these contents, find the corresponding coordinates, and normalize them. After that, we wrap the coordinates in “<bbox></bbox>” and append them to the original wrapped text.

After obtaining the output from GPT-4o, we perform format filtering to remove samples that do not meet the format requirements. And we also perform grounded data checking to correct or remove samples that have incorrect grounding content. We change the grounded blocks “<ocr></ocr><bbox></bbox>” in the questions to “<bbox></bbox>”. Finally, we combine various answers and reasoning to obtain different types of tasks introduced in Sec. 4.

B.2. Instruction Details

As shown in Fig. 9, we introduce the instruction utilized in Multi-granular Parsing tasks and the response format prompts which are followed by the questions for instruction-tuning data. For Bounding Box & Text Localization, we introduce 4 instructions for each task, the answer is directly wrapped with “<ocr></ocr>” or “<bbox></bbox>”. For Full-page Parsing, we introduce 3 instructions for poster data, 1 for chart, and 1 for PDF document data. The responses are in the format shown in Fig. 6 and Fig. 7.

We add different response format prompts to different questions based on the format of the responses to help the model output corresponding results when users interact with it. For generated question and answer pairs, we combine different question and answer pairs with various response format prompts to obtain diverse grounded data. For Grounded Answering Data, we have 7 response format prompts to be added after the question, and the answer should be directly a simple text wrapped with “<ocr></ocr>” and followed by the coordinates wrapped with “<bbox></bbox>”. For Grounded Reasoning Data, we combine two random-chose prompts to-

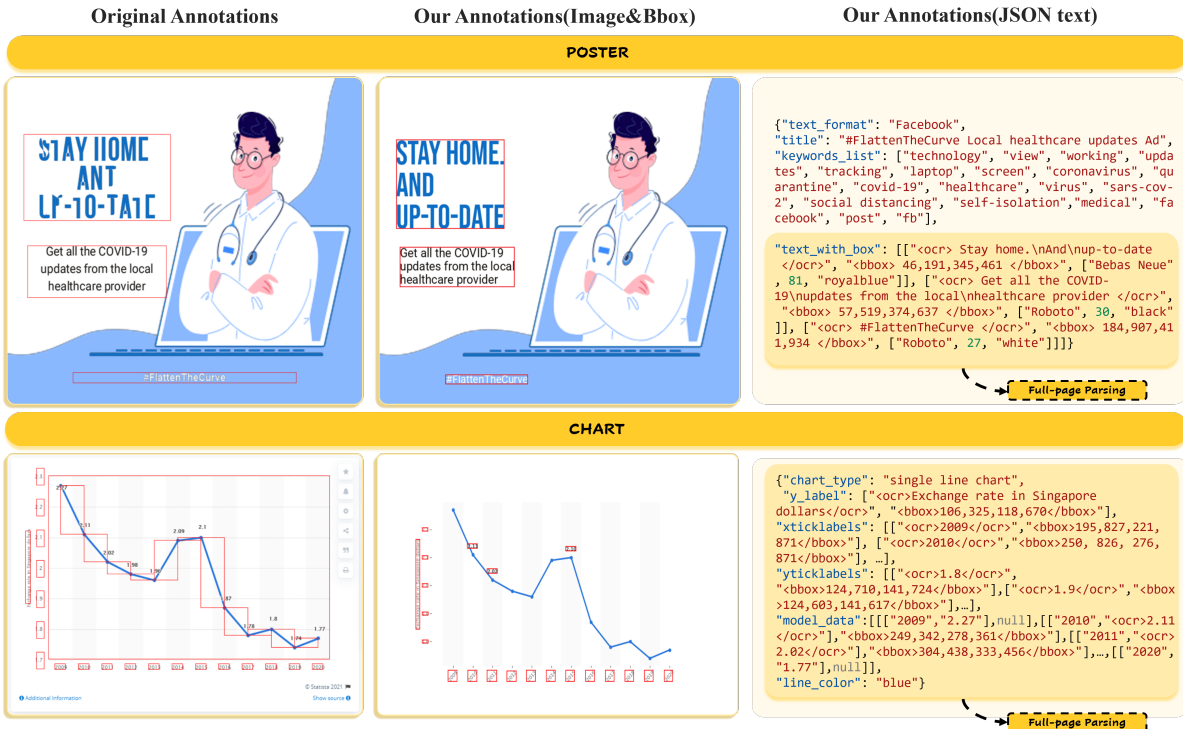


Figure 6. Comparison of origin annotation and our new constructed annotation.

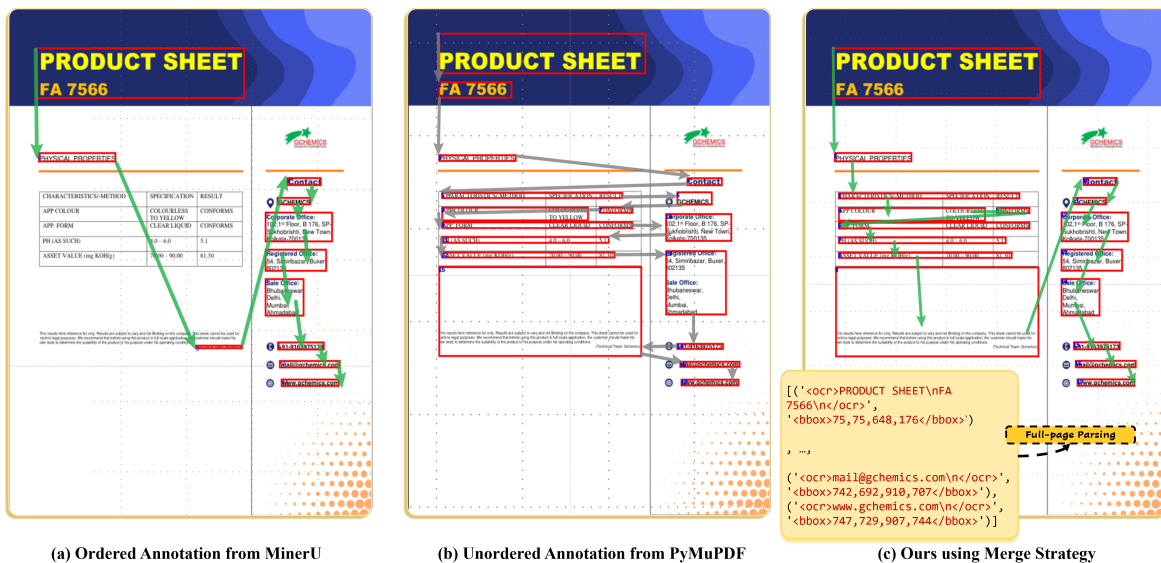


Figure 7. PDF Document Parsing results comparison of ordered annotation from MinerU, unordered annotation from PyMuPDF, and our annotation using the Merge Strategy. Green arrow indicates the ordered annotations and Gray arrow indicates the naive scanning order.

together to form the response format prompt, and the response should be such a sentence structure involves a segment of grounded reasoning followed by "Answer: " and a concise answer. For Grounded Open-ended Answering, we simply use the first part of the response format prompts

for Grounded Reasoning, and the response should be a grounded reasoning sentence. For Plain text Answering, we add no response format prompt to keep consistent with the original question answering.

Poster Q&A Generation Prompts for GPT-4o

System prompt:

You are a poster expert. I'll give you the data about a poster, including some implicit content and a Text box list. Each list item correspond to one box in the poster and represents a segment of text in the poster. The format of each item is [`<ocr>text</ocr>`, `<bbox>bounding box coordinates</bbox>`, (font name, font size, font color)].The list order is somewhat disorganized. You need to reorder these text boxes to make them logically smooth. Please generate 3 most valuable question & answer & explanation & The necessary bbox list to obtain the answer for the content in text box list in json format:

```
[{'question':'','answer':'','explanation':'',' necessary bbox':['<ocr></ocr><bbox></bbox>'}],...].
```

The question types can include factual judgments, factual inquiries, multi-step reasoning, math reasoning, extended understanding, summarize, etc. Note that the implicit content is only for your comprehension for better questions, so don't ask for that content. Note that the response must be related at least one item in the list. In question, you can substitute the text in '`<ocr></ocr>`' with corresponding bbox in '`<bbox></bbox>`', but don't ask too simple questions, like 'What is the text mentioned in `<bbox>A,B,C,D</bbox>`?'.Don't ask about detailed font name or size. When citing a text with corresponding bbox in answer / explanation, enclose it in '`<ocr></ocr>`' followed by its bbox in '`<bbox></bbox>`', formatted as follows: '`<ocr> apple </ocr><bbox>A,B,C,D</bbox>`', where 'A,B,C,D' represents the bbox for 'apple'. The necessary bbox must in your explanation. Assume you directly obtain the poster data based on the poster image, so avoid phrases like 'based on data'. Your answer should be one word or one span, and your explanation should be simple but reasonable.

User prompt:

Here is a poster with the theme [title] in [text_format] style. Here are some related keywords:[keywords_list]. Text box list:[Full-page parsing JSON data]

Chart Q&A Generation Prompts for GPT-4o

System prompt:

You are a chart expert. Here is the data about a chart, including its components and corresponding bounding boxes (bboxes). Additionally, I will give you a question and answer pair to let you know what the chart data represents for. Please generate 3 most valuable question & answer & explanation &The necessary bbox list. Output in JSON format:

```
[{'question':'','answer':'','explanation':'',' necessary bbox':['<ocr></ocr><bbox></bbox>'}],...].
```

The type can include multi step reasoning, math reasoning, extended understanding, color-related reasoning, summarize etc.In question, you can generate pure text question, or question with bbox in '`<bbox></bbox>`' to substitute text in '`<ocr></ocr>`'.

When citing a text or data with corresponding bbox in your answer, enclose it in '`<ocr></ocr>`' followed by its bbox in '`<bbox></bbox>`', formatted as follows: '`<ocr> 50% </ocr><bbox>A,B,C,D</bbox>`' or '`<ocr> 50% </ocr><bbox>null</bbox>`', where 'A,B,C,D' or 'null' represents the bbox for '50%'. The items in "necessary bbox" is the necessary bbox(es) to refer to get the answer, and they must appear in explanation. The items in "necessary bbox" must also be formatted as follows:

'`<ocr> 50% </ocr><bbox>A,B,C,D</bbox>`' or '`<ocr> 50% </ocr><bbox>null</bbox>`', where 'A,B,C,D' or 'null' represents the bbox for '50%'. Assume you directly obtain the data based on the chart image, so avoid phrases like 'based on data', 'according to the dict' or specific color codes (e.g., '#ff3522'). Your answer should be simple.

User prompt: [Full-page parsing JSON data] + [a Q&A sample from ChartQA]

PDF Document Q&A Generation Prompts for GPT-4o

PDF images



+

System prompt:

You are an expert in document reading. Express the original text from the document as much as possible in your reply and use '`<ocr></ocr>`' to enclose the unique origin key words or spans in document.

User Prompt:

Generate a summary of the document and 3 most valuable question & answer & explanation. Make sure your answer is a word or a span, it is correct and explanation is reasonable. The type can include multi-step reasoning, math reasoning, etc. The explanation should be simple but clear.

Note that the output should be in JSON format, like {'summary':'', 'QAs': [{'question':'', 'answer':'', 'explanation':''}, {'question':'', 'answer':'', 'explanation':''}]}

Figure 8. Prompts utilized as input to gpt-4o for poster, chart and PDF document.

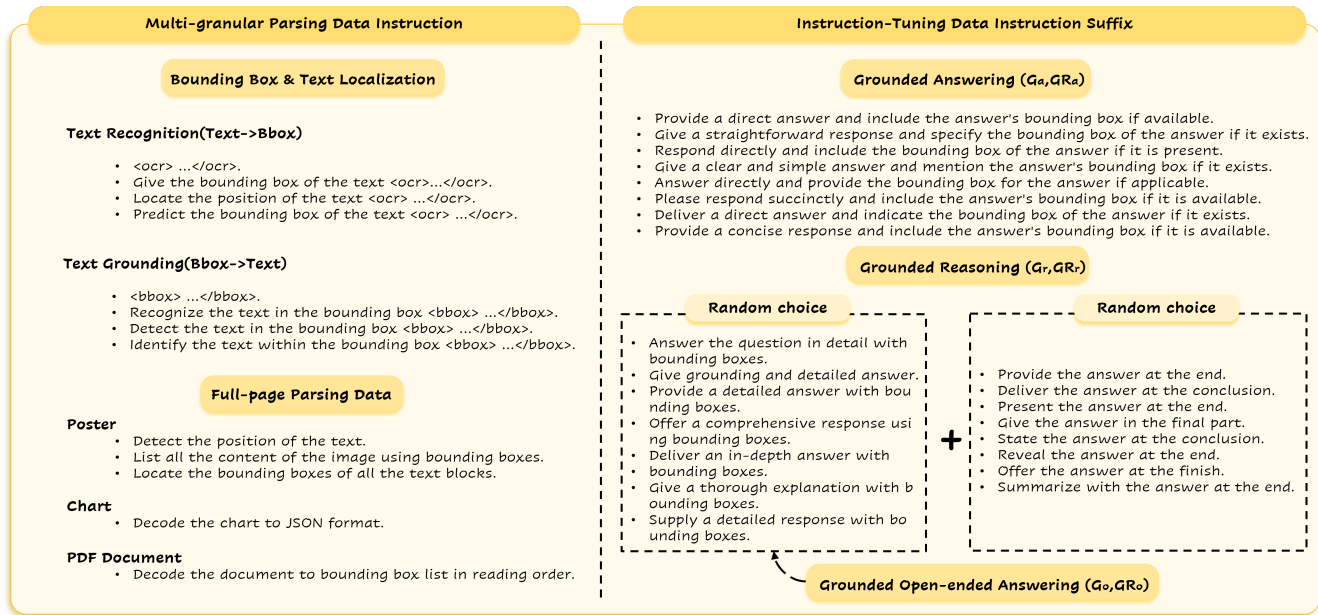


Figure 9. The instruction utilized in Multi-granular Parsing tasks and the response format prompts for instruction-tuning data.

C. Qualitative Results

C.1. Analysis about Fig. 1

We present the inference results of DOGE for grounding, grounding-and-referring, and referring tasks. The specific coordinates of the annotations are omitted, and the grounded text is highlighted with colored boxes. The colors of the boxes within the document image correspond to the colors of the text boxes. For the grounding task, we present four examples. The first sample demonstrates that DOGE can perform fine-grained question answering on general document images with rich content and complex layouts, successfully grounding the corresponding information. The second sample illustrates the model's excellent recognition and localization capabilities for diverse and small text in poster-type images. Additionally, we showcase two samples of grounding in chart-type images, which indicate that DOGE possesses a certain level of mathematical ability, enabling it to provide grounded reasoning during calculation, as well as the capability to estimate values in charts that lack textual annotations based on the axes. For the grounding-and-referring and referring tasks, DOGE is able to recognize the content of user-selected regions and provide reasonable grounded or plain textual reasoning and responses. DOGE exhibits robust fine-grained grounding and referring capabilities, allowing for reliable grounded reasoning and accommodating diverse user interactions, significantly enhancing the overall user experience.

C.2. More Qualitative Results

C.2.1. DOGE-Bench Examples

As shown in Fig. 10, we present three inference examples of DOGE-Bench, demonstrating unlabeled value reading capability, grounding summarization capability, and referring summarization capability. The first figure is an example of reading unlabeled values. We add auxiliary mark boxes in red to the coordinate axes to facilitate reading the values, each box's height represents 5. It can be observed that the model's output values are accurate. The figure on the lower left is an example of grounding summarization for an entire document. DOGE is able to perceive the content of the entire document and perform grounding summarization effectively. The figure on the right is an example of referring summarization for a specific area. The model can also accomplish this task, thereby helping users improve reading efficiency.

C.2.2. Other Examples

It is noteworthy that our model demonstrates strong generalization capabilities beyond some training domains and tasks, as illustrated in Fig. 11.

Strong generalization. As shown in the first row of examples, we feed DOGE with screenshots of the paper content and the specially shaped fan chart that the model does not process before. DOGE is able to correctly respond, demonstrating its strong generalization and usability in actual document reading scenarios.

Handwriting ground-and-refer ability. The middle sample showcases the model's ability to recognize and ground-

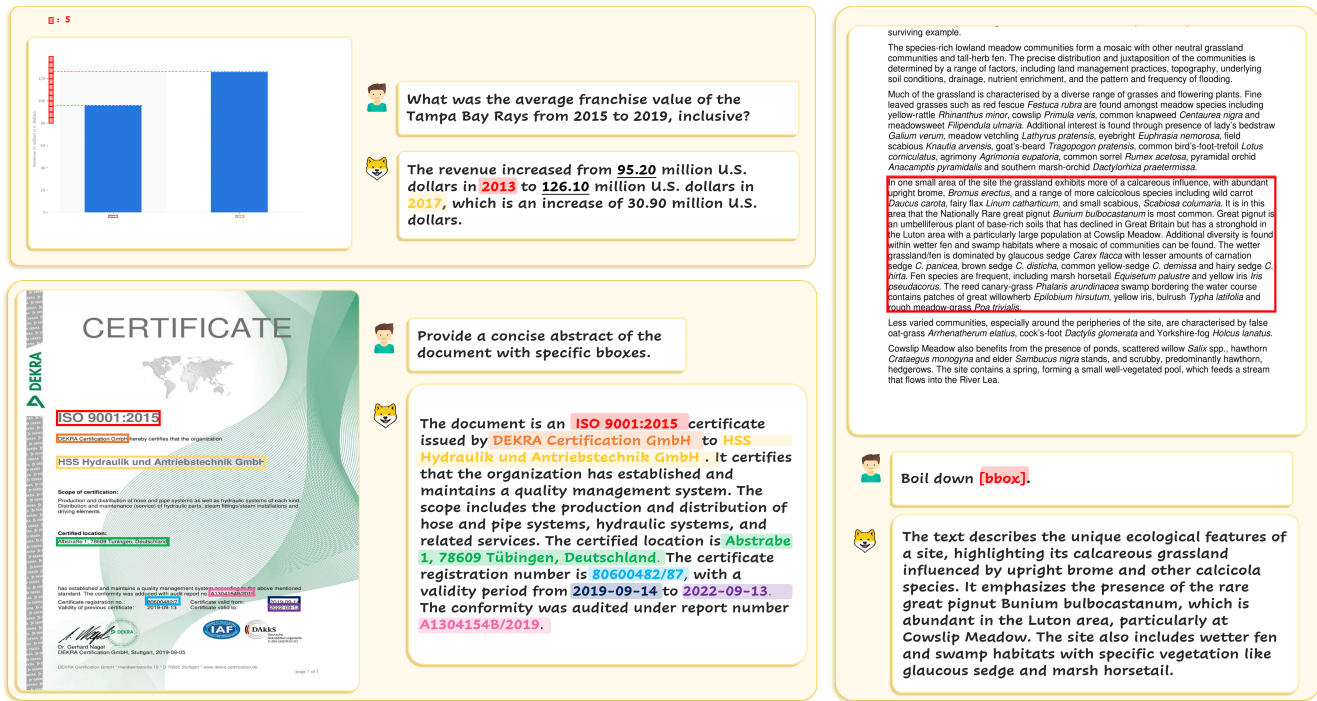


Figure 10. DOGE's inference results on DOGE-Bench.

ing handwritten content. DOGE can fully understand this casual handwriting and provide grounded output. It is worth mentioning that our training data does not contain the grounding or referring tasks on such handwritten images.

Other capabilities. The bottom sample demonstrates an additional untrained capability of our model, specifically referring translation. DOGE is able to provide translations for the asked region. It is important to note that our training data do not include tasks similar to referring translation, even rarely includes languages other than English. Therefore, we believe that DOGE can effectively handle the relationship between regions and corresponding text, and seamlessly integrate the capabilities of large models with grounding and referring abilities.

C.2.3. Failure Cases

Although DOGE demonstrates strong grounding and referring capabilities, there are still some shortcomings, as shown in Fig. 12.

Referring. The upper left figure illustrates an example of incorrect referring. It mistakenly associates the bounding box that should correspond to the text "INITIATIVE" in the question with the text "LOYALTY" below it. Additionally, its width is re-estimated according to the size corresponding to "LOYALTY," resulting in an incorrect answer.

Grounding. When encountering some unfamiliar text content such as tables in the upper right image, DOGE can understand the content, effectively identified the line breaks between "Doc" and "VQA" and merged them together, and

provide correct answers, but the grounding boxes are inaccurate. There is also a issue of incomplete grounding content, such as the bottom sample, which often occurs when the content requiring grounding is interrupted or wrapped to the next line. Although this does not affect understanding, the text and bounding boxes provided by the model do not completely match.

Comprehension. When facing with some unfamiliar structural document, such as the chart shown in the middle right, the model gives incorrect answers. However, due to the intuitive expressiveness of grounding and referring, readers can quickly determine that the model's answer is incorrect. These samples can also provide evidence of the model's deficiencies, laying the foundation for further improvements.

Table 1. Data statistic of DOGE-Bench.

Category	Grounding-and-Referring				Referring		Grounding		Total
	G_1	G_2	G_3	G_4	R_1	R_2	G_5	G_6	
Example	700	800	636	342	427	775	464	4144	

...tinct classes based on both input and output formats. This classification helps in designing clear evaluation metrics. We divide the input formats into two categories based on the presence of bounding boxes: **Grounded Question (GQ)** with bounding boxes, and **Plain-Text Question (PQ)** without bounding boxes. The output formats are categorized into four classes:

- **Grounded Answer (GA)**: The response consists of a brief answer accompanied by its corresponding bounding box.
- **Grounded Reasoning (GR)**: The response includes the detailed reasoning and answer, while the key text contents are grounded.
- **Grounded Open-ended Answer (GO)**: An open-ended response with one or more key text contents grounded, without providing an answer in a certain format.
- **Plain Text Answer (PA)**: This format does not incorporate grounded text content.

By combining two input forms and four output forms, we derive seven document referring and grounding tasks. Among these tasks, three sub-tasks primarily assess grounding capability: grounded answering for plain-text questions (G_1), grounded reasoning for plain-text questions (G_2), and grounded open-ended answering for plain-text questions (G_3). The plain-text answering for grounded questions (G_4), and grounded reasoning for grounded questions (G_5), and grounded open-ended answering for grounded questions (G_6), require the integration of both grounding and referring capabilities for successful completion.

Metrics. Our benchmark evaluation encompasses two aspects: grounding performance and text answer accuracy. Following the previous works in grounded captioning [47], we evaluate the grounding and text answer separately. For grounding performance, we use $F_{1\mu}$ score, which evaluates grounding results as a multi-label classification problem. The generated grounded text is considered correct if Intersection over Union (IoU) between its bounding box and the GT bounding box is greater than 0.5, meanwhile its text matches with the GT text. For text answer accuracy, we use exact text matching accuracy for short-answer tasks and BLEU score for long-answer tasks.

Data Statistics. Our DOGE-Bench includes 2K grounding samples, 0.5K referring samples, and 15K grounding-and-referring samples. The detailed statistics is shown in Tab. 1.

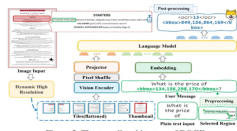
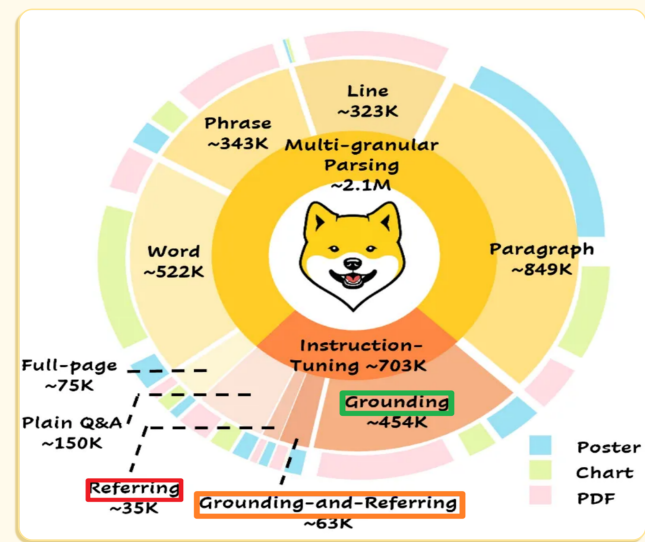


Figure 5. The overall architecture of DOGE.

5. DOGE
Overall Architecture. As illustrated in Fig. 5, our model DOGE employ a general MLLM architecture, including a vision encoder, a projector, and a large language model. In the vision encoder component, to enable the model to handle high resolution, we first search for the best aspect ratio for the input image and dynamically segment images into multiple tiles. These tiles, along with a thumbnail of the input image, are provided as input to the vision encoder. We employ pixel shuffle [9] to improve the computational efficiency of the model when processing high-resolution images. For bounding box representation, we simply discretize the continuous coordinates into discrete values from 0 to 999, avoiding the introduction of extra module or location tokens. During inference, we transfer the user-selected regions to coordinates of bounding boxes and insert them into the query for preprocessing. After obtaining the output, we utilize post-processing to overlay bounding boxes on original document images, thereby facilitating user interaction.

Training Strategy. We adopt a three-stage training strategy, including pre-aligning, pre-training, and fine-tuning. The pre-aligning stage focuses on aligning the feature space of vision and language rapidly. In this stage, we freeze both vision encoder and large language model, and train the projector with a relatively large learning rate. The pre-training stage aims at document parsing capabilities. We unfreeze the vision encoder and the LLM, enabling the model to recognize diverse textual content and acquire text-reading capability. In the fine-tuning stage, we train the entire model using diverse instruction-tuning data, enhancing its instruction-following ability while activating its grounding and referring capabilities during dialogue and reasoning.

Training Dataset. In the pre-aligning stage, we utilize LLaVA-558K [21] to train the projector. For the pre-training dataset, we utilize DocStruct4M [10] along with our 1.4M multi-granular parsing data to enhance basic text reading, text grounding and text referring capabilities of DOGE. For the fine-tuning data, we adopt our ground-and-refer instruction-tuning data and meticulously selected



How many output formats are there?

There are four output formats:
Grounded Answer (GA),
Grounded Reasoning (GR),
Grounded Open-ended Answer (GO),
and
Plain Text Answer (PA).

What's the sum of [bbox], [bbox], and [bbox]?

The sum of Referring, Grounding-and-Referring, and Grounding is 35K + 63K + 454K = 552K.

Q & A:

1. Is DOGE a powerful document understanding assistant?
 Yes! It has Grounding and Referring capabilities, which can help readers read efficiently.

2. Can DOGE handle handwritten content?
 Yes! It can, as you can see right now.

Why is DOGE useful?

DOGE is useful because it has Grounding and Referring capabilities, which help readers read efficiently.

Can DOGE recognize manuscript?

Yes, DOGE can handle handwritten content. This is confirmed by the text 'Yes! It can, as you can see right now.'.

1. Introduction

In recent years, Multimodal Large Language Models (MLLMs) [5, 17, 14, 15, 20, 31, 46–48] have achieved significant advancements in general visual understanding and reasoning by integrating pre-trained vision encoders with large language models. Leveraging the diverse fine-grained annotations of existing image datasets, some researchers further equip MLLMs with grounding and referring capabilities. These capabilities enhance the detailed visual understanding, increase the credibility of responses, and facilitate more efficient human-AI interaction.

In the field of visual document understanding, the dense textual content and complex layout significantly complicate fine-grained understanding. To build a user-friendly and trustworthy document AI assistant, it is crucial to enable users to refer to specific regions of a document for more precise comprehension and provide accurate grounding of key details for more expressive and effective interaction.

However, due to the lack of high-quality, fine-grained document datasets, the full potential of MLLMs in document grounding and referring remains largely unexplored. Existing efforts leverage the multi-granularity parsing annotations to enhance document detailed perception [10], or introduce region-level instruction-tuning tasks to achieve basic referring capabilities [18]. However, these works have no significant shortcomings.

Translate [bbox] into Chinese.

然而，由于缺乏高质量的细粒度文档数据集，MLLMs在文档定位和引用方面的潜力仍然未被充分利用。现有努力利用多粒度解析注释来增强文档详细感知[10]，或引入区域级指令调谐任务以实现基本引用能力[18]。然而，这些工作有两个显著的局限性：

Figure 11. Other DOGE's inference samples.

Incorrect Referring

CHARACTER CORE

ALERTNESS
Being aware of what is taking place around me so I can respond appropriately

ATTENTIVENESS
Concentrating on the person or task before me

AVAILABILITY
Willingness to change my schedule and priorities to meet a need

CAUTIOUSNESS
Taking time to ensure the right decision is made or action is taken

COMPASSION
Helping those in need

COOPERATION
Understanding others so I can effectively work with them

COURAGE
Overcoming fear by saying and doing what is right

CREATIVITY
Approaching a need, a task, or an idea from a new perspective

DECISIVENESS
Processing information and finalizing difficult decisions

DEPENDABILITY
Fulfilling commitments even in the face of difficulty

DETERMINATION
Overcoming obstacles in order to reach my goal

DILIGENCE
Focusing my effort on the work at hand

DISCIPLINE
Choosing behaviors to help me reach my goals

ENDURANCE
The inner strength to withstand stress and do my best

ENTHUSIASM
Expressing interest and excitement in what I do

FLEXIBILITY
Adjusting to change with a good attitude

FORGIVENESS
Releasing feelings of resentment

GENEROSITY
Managing resources to freely give

GRATEFULNESS
Demonstrating appreciation to others for what I have and how they have helped me

HONESTY
Being truthful in what I say and do

HUMILITY
Recognizing the people and factors that have shaped my life

INITIATIVITY
Recognizing and doing what needs to be done before I am asked to do it

LOYALTY
Demonstrating commitment to others

ORDERLINESS
Organizing my thoughts and surroundings for greater achievement

PATIENCE
Taking the time necessary to work through a difficult situation

PERSUASIVENESS
Effectively communicating with others so they can better understand

POSITIVITY
Maintaining a good attitude, even when faced with difficulty

PUNCTUALITY
Showing respect for others by doing the right thing at the right time

RESILIENCE
Recovering from adversity

RESPECT
Treating others with honor and dignity

RESPONSIBILITY
Knowing and doing what is expected of me

SINCERITY
Doing what is right with transparent motives

THOROUGHNESS
Taking care of necessary details

TOLERANCE
Demonstrating respect for others who do not share my perspective

TRUSTWORTHINESS
Gaining the confidence of others by demonstrating reliability

WISDOM
Making practical application of what is learned

ENGAGE YOUR TEAM. INSPIRE YOUR PEOPLE. CREATE A WINNING CULTURE.
© 2015 Strata Leadership LLC • StrataLeadership.com • 877.387.0001



How is [bbox] described in the document?



LOYALTY is defined as the commitment to others, emphasizing the importance of being dedicated and supportive in relationships.

Answer:

Demonstrating commitment to others

Correct Answer:

Recognizing and doing what needs to be done before I am asked to do it.

Inaccurate Grounding

Table 3. Performance comparison on 10 general document benchmarks.

Model	Size	DocVQA	InfoVQA	Deep Form	KLC	WTQ	Tab Fact	Chart QA	Text VQA	Text Caps	Visual MRC
Domu[13]	1B	67.5	11.6	61.6	30.0	18.8	54.6	41.8	43.5	74.4	93.91
DocOwl[9]	7B	62.2	38.2	42.6	30.3	26.9	60.2	57.4	52.6	111.9	188.8
URReader[43]	7B	65.4	42.2	49.5	32.8	29.4	67.6	59.3	57.6	118.4	221.7
TextMonkey[16]	9B	73.0	28.6	59.7	37.8	31.9	-	66.9	65.9	-	-
Vary[39]	7B	76.3	-	-	-	-	-	66.1	-	-	-
TokenPacker[40]	13B	70.0	-	-	-	-	-	-	-	-	-
DocPeida[8]	7B	47.1	15.2	-	-	-	-	46.9	60.2	-	-
QwenVL[1]	9B	65.1	35.4	-	-	-	-	65.7	63.8	-	-
IXC 2.5[50]	7B	90.9	70.0	71.2	-	53.6	85.2	82.2	78.2	-	307.5
DocOwl-1.5-Chat[10]	8B	82.2	50.7	68.8	38.7	40.6	80.2	70.2	68.6	131.6	246.4
DocOwl-2[11]	8B	80.7	46.4	66.8	37.5	36.5	78.2	70.0	66.7	131.8	217.4
InternVL2[6]	8B	91.6	74.8	-	-	-	-	83.3	77.4	-	-
DOGE	8B	91.7	70.7	70.8	40.4	58.8	84.5	83.6	76.6	145.9	332.5

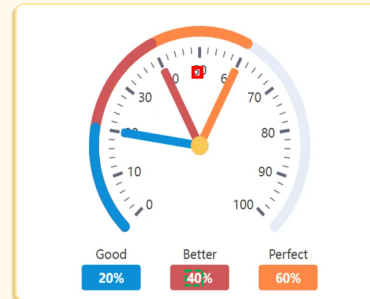


Which method has the best performance on DocVQA?



The method with the best performance on DocVQA is **DOGE**, achieving a score of 91.7.

Incorrect comprehension



What's the value of the red pointer?



The red pointer is at 50.

Correct Answer:

The red pointer is at 40.

Incomplete Grounding

CICR

Search results for "DROIT INTERNATIONAL PENAL" (Keyword in French)

Document type:	Periodical Issue
Title:	Punishment and society : the international journal of penology [Vol. 14, no. 5 December 2012]
Languages:	English
Keyword in English:	DEPRIVATION OF LIBERTY CRIMINAL LAW INTERNATIONAL CRIMINAL LAW
Keyword in French:	PRIVATION DE LIBERTE DROIT PENAL DROIT INTERNATIONAL PENAL
Go to:	Punishment and society : the international journal of penology

Issue text	Issue date	Status	Item number	Call number / Mark	Location	Disposability	Due date
Vol. 14, no. 5, December 2012	01.02.2013	received	100032202	PER	BIB	Available	



What is the title of the periodical issue related to 'DROIT INTERNATIONAL PENAL'?



Punishment and society : the international journal of penology [Vol. 14, no. 5, December 2012]

Correct Answer:

Punishment and society : the international journal of penology [Vol. 14, no. 5, December 2012]

Figure 12. Failure cases of DOGE.

Table 6. The detailed composition of the DOGE dataset’s Multi-granular parsing data and Instruction-Tuning data.

Type	Number	SubType	Number	SubSubType	Number
Multi-granular Parsing	2,114,414	Full Page	75,391	Poster	20,867
				Chart	31,716
				PDF	22,808
		Word	522,682	Poster	65,446
				Chart	354,731
PDF	102,505				
Span	343,596	Poster	56,006		
		Chart	58,212		
PDF	229,378				
Line	323,211	Poster	5,577		
		Chart	6,268		
PDF	311,366				
Paragraph	849,534	Poster	511,998		
		Chart	229,378		
		PDF	108,158		
Instruction-Tuning	703,724	Grounding	454,404	Poster	96,718
				Chart	62,636
				PDF	295,050
		Grounding + Referring	63,243	Poster	36,663
				Chart	849
PDF	25,731				
Referring	35,197	Poster	19,207		
		Chart	618		
		PDF	15,372		
Plain Text Q&A	150,880	Poster	35,831		
		Chart	45,948		
		PDF	69,101		

D. Dataset Details

D.1. Dataset Statistic

As shown in Tab. 6, we present the detailed composition of the DOGE dataset’s Multi-granular Parsing data and Instruction Tuning data. The Multi-granular Parsing data includes five granularities: word, phrase, line, paragraph, and full-page parsing. The Instruction-Tuning data comprises three types of grounded data: grounding referring, grounding-and-referring, and plain text Q&A. We provide the detailed data volume for each type.

Fig. 13 shows the distribution of block counts for poster, chart, and PDF document. Poster has fewer grounded blocks within each image, larger areas, diverse font styles, and larger font sizes. The average text length per block for poster is 3.25 words. Chart, on the other hand, has a higher number of blocks with small areas and small font sizes. Each block in a chart corresponds to a small component value within the chart, with an average text length of 1.34 words per block. PDF document has a moderate distribution of block counts, larger areas, and small font sizes. Each block in a PDF contains relatively longer text, with an average length of 22.55 words per block.

D.2. Construction Cost Details

In terms of multi-granular parsing data construction, DOGE-Engine can achieve boundary box annotation and content extraction on any document dataset with a similar rendering process. It can also scale up to a larger volume of PDF source files without additional manual costs. When constructing instruction fine-tuning data, our main cost lies

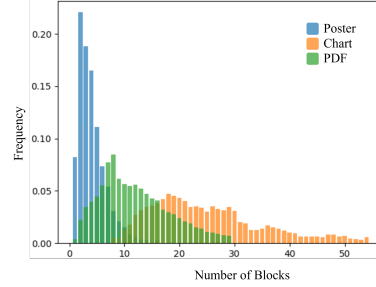


Figure 13. Block count distribution of poster, chart, and PDF.

Table 7. The detailed statistics of 716k other document-related data.

Dataset	# Samples	Dataset	# Samples
IIT5K [41]	1,990	RoBUT WTQ[70]	38,241
TextOCR-GPT4V [3]	25,104	AI2D (InternVL [6])	12,403
FigureQA [16]	1,000	Infographic VQA [40]	8,489
Diagram Image2Text	295	LRV Chart [30]	1,776
K12 Printing	20,000	SROIE	33,616
AI2D (GPT4V Detailed Caption) [17]	4,864	MultiHiertt	7,614
VisText [51]	9,964	RoBUT WikiSQL	74,984
ChartQA [37]	18,260	VisualMRC[50]	3,022
DVQA [15]	20,000	TextCaps [46]	21,942
Magpie Pro [58] (L3 ST)	50,000	Chart2Text [42]	26,956
HiTab [8]	2,495	HME100K [66]	74,492
RoBUT SQA	8,509	Magpie Pro (L3 MT)	50,000
ChromeWriting [57]	8,825	Magpie Pro (Qwen2 ST)	50,000
Screen2Words [54]	15,725	Rendered Text [57]	9,995
IAM [36]	5,658	TQA [18]	27,302
AI2D (Original)	2,429	SynthDog-EN [20]	40,000
MMC bInstruction Arxiv[31]	20,000	MMC bInstruction NON-Arxiv	20,000

in the API calls to GPT-4o. We use the version gpt-4o-2024-08-06. For poster and chart, our input is long text containing full page parsing data, while for PDF, the input is document images plus short prompt texts. During construction, we can generate multiple question and answer pairs for a single image simultaneously and use batch API requests to save costs. Ultimately, we can construct over 1,000 grounded question and answer pairs for an average cost of \$1, which is far more efficient and cost-effective than manual construction.

E. Training Data Details

For the pre-training data, we utilize DocStruct4M [12] along with our 2.1M multi-granular parsing data to enhance DOGE’s foundational grounding and referring capabilities.

For the fine-tuning data, we compose our 703k Instruction-Tuning data with 575k DocDownStream-1.0 data from [12] and 716k other document-related data from various datasets, resulting in a final fine-tuning dataset of 2M. In Tab. 7, we show the detailed data source and sampled number of 716k other document-related data.

F. More Experiments

Effectiveness of our pretraining data. To better assess base performance of pretrained model using our multi-

Table 8. Grounding and recognition performance comparison of the same model pretrained w. and w.o. our pretraining data.

	Poster					PDF Document					Chart
Text Grounding											
Model	word	phrase	line	paragraph	ALL	word	phrase	line	paragraph	ALL	ALL
DOGE ⁻	58.0	63.12	58.96	76.43	63.19	33.88	47.62	51.5	62.88	48.97	24.03
DOGE	91.88	94.62	96.28	91.04	93.53	70.25	82.25	87.5	83.62	80.91	86.12
Text Recognition											
DOGE ⁻	83.34	60.96	41.67	36.95	55.73	49.66	58.62	50.34	47.27	51.47	68.86
DOGE	92.94	94.05	86.52	86.08	89.9	73.57	86.48	76.85	71.8	77.17	95.24

granular parsing data, we further introduce our text&bbox localization test set **DOGE-Local15k** comprising four granularities (word, phrase, line, paragraph) across three categories of data: poster, chart, and general document, similar to DocLocal4k[12]. DOGE⁻, same as Tab. 4, is pre-trained without using our constructed multi-granular parsing data. For the chart data, we claim that the chart type annotations in DocStruct4M for during DOGE⁻'s pre-training are inconsistent with our grounding training objectives, which aims to match the text and the bbox. In the original data, the text corresponds to bounding boxes of bar/line charts. Therefore, these value is not referenceable, and we mark them in gray. Moreover, in the case where the granularity of data categories for charts is relatively singular, we directly reported the results for "ALL".

As shown in Tab. 8, after incorporating our multi-granular parsing data for pre-training, the model exhibits enhancements in text recognition and grounding tasks. Our improvements on DOGE-Local15k are significant, thereby laying a solid foundation for accurate document grounding interactions.