MMDocBench: BENCHMARKING LARGE VISION LANGUAGE MODELS FOR FINE-GRAINED VISUAL DOC UMENT UNDERSTANDING

Anonymous authors

Paper under double-blind review



Figure 1: Overview of MMDocBench, which is designed to holistically assess the fine-grained visual understanding capability of LVLMs through various document understanding tasks. It consists of 15 main tasks and 48 sub-tasks, involving 2,400 document images, 4,338 QA pairs and 11,353 supporting regions (*i.e.*, bounding boxes). Each QA pair corresponds to one or more supporting regions, marked with red dotted-line rectangles on the images.

ABSTRACT

Large Vision-Language Models (LVLMs) have achieved remarkable performance in many vision-language tasks, yet their capabilities in fine-grained visual understanding remain insufficiently evaluated. Existing benchmarks either contain limited fine-grained evaluation samples that are mixed with other data, or are confined to object-level assessments in natural images. To holistically assess LVLMs' fine-grained visual understanding capabilities, we propose using document images with multi-granularity and multi-modal information to supplement natural images. In this light, we construct MMDocBench, a benchmark with various OCR-free document understanding tasks for the evaluation of fine-grained visual perception and reasoning abilities. MMDocBench defines 15 main tasks with 4,338 QA pairs and 11,353 supporting regions, covering various document images such as research papers, receipts, financial reports, Wikipedia tables, charts, and infographics. Based on MMDocBench, we conduct extensive experiments using 10 open-source and 3 proprietary advanced LVLMs, assessing their strengths and weaknesses across different tasks and document image types. The benchmark, task instructions, and evaluation code will be made publicly available.

1 INTRODUCTION

055 056

057 Large Vision-Language Models (LVLMs) have attained remarkable performance across various 058 vision-language tasks (Yin et al., 2024). However, existing LVLMs, such as GPT-4V (OpenAI, 2023), LLaVA (Liu et al., 2023a) and MiniGPT-4 (Zhu et al., 2023), still struggle with understanding fine-grained visual details in images. For instance, (Tong et al., 2024) have demonstrated that 060 many LVLMs perform poorly in visual grounding to image details for visual question answering. 061 This fine-grained visual understanding capability is indispensable for LVLMs in many downstream 062 tasks (Peng et al., 2024; Xuan et al., 2024), such as object recognition (Lin et al., 2014), image 063 segmentation (Minaee et al., 2022) and forgery detection (Qu et al., 2023). As such, it is essential to 064 develop LVLMs' capability in fine-grained visual understanding. 065

To achieve this goal, a key prerequisite is establishing benchmarks that can comprehensively eval-066 uate the strengths and weaknesses of LVLMs in fine-grained visual understanding. However, rep-067 resentative multimodal benchmarks such as MMVet (Yu et al., 2023), MME (Fu et al., 2024), and 068 MMT-Bench (Ying et al., 2024) contain relatively few data samples to examine LVLMs' under-069 standing of fine-grained details rather than the entire image. Besides, these samples are not isolated from the overall dataset. This makes it difficult to evaluate LVLMs' ability in fine-grained visual 071 understanding comprehensively. Moreover, some benchmarks, such as Visual7W (Zhu et al., 2016), 072 RefCOCO (Yu et al., 2016), GVT-bench (Wang et al., 2023a), and MMBench (Liu et al., 2023b), 073 have designed specific tasks like object grounding and object counting to evaluate fine-grained vi-074 sual understanding. However, these tasks are confined to object-level details in natural images and 075 do not assess finer-grained details.

076 To address the limitations, we consider using document images to supplement natural images for 077 the evaluation of fine-grained visual understanding. As illustrated in Figure 1, document images encapsulate various document content types, including text, figures, tables, charts, and diagrams, 079 all presented within a visual format. These document images like receipts and research papers are widely used across various domains such as finance, legal, education, and academia (Kim et al., 081 2022). Compared to natural images, document images offer certain advantages as testing data to evaluate LVLMs' fine-grained visual understanding capabilities. In particular, 1) document images 083 contain different granularities of information to evaluate the *fine-grained visual perception* abilities, such as the localization and recognition of text and tables. The diverse elements (e.g., text, table, and 084 chart) tend to occupy only a small part of the entire image, yet convey critical information in various 085 granularities; for example, a receipt image may comprise item description (sentence level), quantity (token level), and barcode (object level). Moreover, 2) document images require LVLMs to integrate 087 multi-granularity and multi-type information to perform complex reasoning, thereby evaluating their 088 fine-grained visual reasoning abilities. For example, in the bar chart located in the bottom left of 089 Figure 1, LVLMs need to combine the legend, axis labels, and numerical values on the bar chart for comprehensive reasoning. 091

In light of this, we construct a benchmark, MMDocBench, using document images to assess LVLMs' 092 capabilities in fine-grained visual document understanding. We design various OCR-free document understanding tasks from the perspectives of *fine-grained visual perception* and *fine-grained visual* 094 reasoning, where understanding partial and fine-grained details in images is crucial, rather than treat-095 ing the image as a whole. For fine-grained visual perception, MMDocBench encompasses nine tasks 096 to evaluate the LVLMs' capabilities in fine-grained information recognition, localization, detection, 097 and extraction, including Text Recognition, Table Recognition, Text Localization, Table Cell Local-098 ization, Key Information Extraction, Document Forgery Detection, Document Question Answering (QA), Chart QA, and Infographic QA. For fine-grained visual reasoning, we design six tasks to assess the reasoning ability by integrating fine-grained information: Arithmetic Reasoning, Logical 100 Reasoning, Spatial Reasoning, Counting, Comparison, and Sorting. We present one example for 101 each task in Figure 1 and refer to Section A.4 in Appendix for more comprehensive examples. 102

Based on these tasks, we select document images from 21 document understanding datasets to con struct QA pairs for evaluation. The document images span a wide variety of types, including re search papers, book covers, financial reports, scene-text images, receipts, Wikipedia tables, charts,
 infographics, and other industry documents. Additionally, in MMDocBench, we also provide anno tations of supporting regions (*i.e.*, bounding boxes) within the image for each QA pair, as shown by
 red dotted rectangles in Figure 1. The supporting regions enable the evaluation of whether LVLMs

have correctly grounded their predictions on the associated regions in the image, leading to a more comprehensive evaluation. Furthermore, the output of supporting regions offers significant practical value, making the LVLMs' responses more informative and interpretable while allowing for rapid cross-checking between the answer and the image. Finally, MMDocBench contains 2,400 document images, involving 4, 338 QA pairs with 11, 353 supporting regions.

113 With MMDocBench, we evaluate 10 open-source and 3 proprietary LVLMs that have demonstrated 114 impressive performance in many vision-language tasks. Experimental results reveal that our MM-115 DocBench poses significant challenges to current LVLMs, especially in region prediction, where 116 almost all LVLMs struggle to identify the supporting regions. Although a notable gap persists in 117 answer prediction between open-source and closed-source LVLMs, their performance in region pre-118 diction is comparable. The in-depth analyses on LVLMs' performance across different tasks and document types provide more insights: 1) Localization and detection tasks present significantly 119 greater challenges compared to other tasks; 2) Region prediction on document images containing 120 diverse elements (e.g., table, chart, and figure) is typically more challenging than on general docu-121 ments, whereas answer prediction shows no significant variation across different document types. 122

- 122
- 123 In summary, our contributions are summarized as follows.
- We highlight the gap in existing benchmarks for evaluating LVLMs' fine-grained visual understanding capabilities. Besides, we propose leveraging document images with multi-granularity and multi-type information to complement natural images in assessing the fine-grained visual perception and reasoning abilities of LVLMs.
 - We construct MMDocBench, a comprehensive benchmark for tracking LVLMs' progress in finegrained visual document understanding. MMDocBench defines 15 main tasks and 48 sub-tasks over a wide range of document types, involving 2, 400 document images, 4, 338 QA pairs, and 11, 353 supporting regions for a holistic evaluation.
 - We conduct extensive experiments with 13 representative LVLMs on MMDocBench. We report the major strengths and weaknesses of LVLMs in different tasks and document types and provide valuable insights, facilitating the advancements of LVLMs in future.
- 135 136 137

138 139

140

129

130

131

132

133

134

2 RELATED WORK

2.1 LARGE VISION-LANGUAGE MODELS

141 Large Vision Language Models (LVLMs) are built upon Large Language Models (LLMs) like 142 GPTs (Radford et al., 2019; Brown et al., 2020) and LLaMA (Touvron et al., 2023). A common approach is to utilize a visual encoder to encode the image first and then apply a visual adapter to 143 align visual and textual representations within the LLMs. LVLMs are typically trained in two stages, 144 i.e. self-supervised pre-training over large-scale image-text pairs and supervised instruction tuning 145 with annotated data. Notable examples include Flamingo (Alayrac et al., 2022), BLIP-2 (Li et al., 146 2023a), LLaVA (Liu et al., 2023a), MiniGPT-4 (Zhu et al., 2023), mPLUG-Owl (Ye et al., 2023c), 147 CogVLM (Wang et al., 2023b), MiniCPM-V (Yao et al., 2024), Monkey (Li et al., 2024b), and 148 InternVL (Chen et al., 2024d). For instance, BLIP-2 (Li et al., 2023a) utilizes frozen CLIP (Rad-149 ford et al., 2021) as its visual encoder and proposes a Q-Former as the visual adapter. LLaVA (Li 150 et al., 2024a) and MiniGPT-4 (Zhu et al., 2023) replace the visual adapter with simple linear layers, 151 demonstrating impressive effectiveness.

152 Recently, increasing efforts have been devoted to enhancing LVLMs' performance in document im-153 age understanding (Ye et al., 2023b; Liu et al., 2024b; Ye et al., 2023a; Liu et al., 2024c; Bai et al., 154 2023; Zhang et al., 2023), generally by accommodating high-resolution input images or/and improv-155 ing quality of training corpus. For instance, LLaVA-NeXT (Liu et al., 2024b) extends LLaVa (Liu 156 et al., 2023a) to enhance its OCR capability by increasing image resolution and mixing more docu-157 ment images in visual instruction tuning data; TextMonkey (Liu et al., 2024c) and MiniCPM-V (Yao 158 et al., 2024) divide high-resolution images into window patches with a sliding window method and 159 reduce the length of visual tokens via token compression techniques. Compared with natural images, the understanding of document images requires LVLMs' interpretation of fine-grained image 160 details. In this work we propose a comprehensive benchmark to evaluate such fine-grained visual 161 understanding capability of LVLMs to facilitate future research.



Figure 2: An illustration of benchmark construction pipeline for MMDocBench.

2.2 EVALUATION BENCHMARKS FOR LVLMS

173 To date, lots of multimodal benchmarks have been constructed to holistically assess LVLMs' inte-174 grated capabilities, such as recognition, knowledge, math, reasoning and safety (Huang & Zhang, 175 2024). Some representative examples include LVLM-eHub (Xu et al., 2023), MME (Fu et al., 176 2024), MMStar (Chen et al., 2024a), SEED-Bench (Li et al., 2023b), MMMU (Yue et al., 2024a), 177 MMMU-Pro (Yue et al., 2024b), MMVet (Yu et al., 2023) and MMT-Bench (Ying et al., 2024), etc. 178 In addition, there are also some benchmarks designed to evaluate one specific aspect of LVLMs. 179 For instance, POPE (Li et al., 2023b) and HallusionBench (Guan et al., 2024) focus on evaluating hallucination; MathVista (Lu et al., 2024) is centered on assessing the visual mathematical reason-181 ing capability; OCRBench (Liu et al., 2023c) assesses the OCR capability on document images 182 with five text-related visual tasks. These existing benchmarks include some samples that can be used to examine LVLMs' fine-grained visual understanding capability, but they are not only limited 183 in number, but also difficult to be separated from other samples. Moreover, their evaluations rely 184 solely on natural language output without supporting region prediction involved, which is a crucial 185 measure for achieving fine-grained visual understanding in LVLMs (Peng et al., 2024). Current available benchmarks specially for evaluating fine-graded visual understanding of LVLMs mainly 187 include Visual7W (Zhu et al., 2016), RefCOCO (Yu et al., 2016), GVT-bench (Wang et al., 2023a), 188 and MMBench (Liu et al., 2023b), which are however centered on object-level recognition and in-189 terpretation within natural images, offering limited information granularity. Our MMDocBench is 190 the first comprehensive benchmark aiming to evaluate LVLMs' fine-grained visual understanding 191 capability with various OCR-free document understanding tasks.

192 193

194

169

170 171 172

2.3 VISUAL GROUNDING IN LVLMS

Visual Grounding (VG) is a technique that localizes the relevant object or region to a given natural 195 language description in the visual input (Deng et al., 2018). It has been applied in LVLMs to fa-196 cilitate generating more informative and comprehensive responses, benefiting various downstream 197 applications like Object Recognition (Lin et al., 2014), Image Segmentation (Minaee et al., 2022) and Referential Comprehension (Yu et al., 2016). Existing grounded LVLMs can be classified into 199 two categories: 1) region-level (i.e., bounding box) grounding based LVLMs, such as OFA (Wang 200 et al., 2022), Kosmos-2 (Peng et al., 2024) and Pink (Xuan et al., 2024); and 2) pixel-level grounding 201 based LVLMs, such as PixelLM (Ren et al., 2024), GlaMM (Rasheed et al., 2024), Lisa (Lai et al., 202 2024) and Ferret (You et al., 2024). However, the training and evaluation of these LVLMs sorely rely 203 on object-level tasks involving natural images, overlooking the various granularities in document 204 images. In this work, we build MMDocBench to foster the advancement of the fine-grained visual 205 understanding capability in LVLMs with a variety of OCR-free document understanding tasks. In 206 MMDocBench, it requires LVLMs to perform region-level grounding following typical settings in most document understanding tasks, such as Optical Character Recognition (OCR) (Singh et al., 207 2021), Key Information Extraction (Huang et al., 2019) and Table Recognition (Zheng et al., 2021). 208

209 210

3 PROPOSED MMDOCBENCH

211 212

213

3.1 PROBLEM DEFINITION

214 On MMDocBench, the LVLMs are expected to output the precise answer to a natural language query 215 given a document image while highlighting the supporting regions within the image contributing to inferring the answer. Formally, given a document image *d* possibly containing text, table, chart, and/or figure, for a question q, a Large Vision-Language Model \mathcal{M} is required to produce a response including the answer a and the corresponding regions R as supporting evidence, formulated as

$$\mathcal{M}(d,q) = (a,R). \tag{1}$$

Each region in R is a bounding box that is represented by the coordinates of its top-left and bottomright corners in the format of $[x_1, y_1, x_2, y_2]$.

3.2 CONSTRUCTION PIPELINE

219 220

221

222 223

224 225

226

227

228 229

230

231

232

233

234

235

The pipeline for constructing our MMDocBench is illustrated in Figure 2.

Step 1: Taxonomy Design. Targeting at a comprehensive evaluation of LVLMs' fine-grained visual understanding capabilities, we design the taxonomy of MMDocBench following two principles.

- **Fine-grained Discrimination**: The MMDocBench must provide tasks that can adeptly evaluate LVLMs' visual comprehension capabilities with sufficient discriminability of fine-grained details in the image, rather than treating the image as a whole.
- **Diversity**: The MMDocBench must encompass a broad range of tasks in terms of required capabilities (e.g., perception, reasoning), content granularity (e.g., characters, words, tables), and document types (e.g., scientific papers, financial reports, receipts).

236 To solve the problem in MMDocBench, two capabilities are at the core for LVLMs, i.e. fine-grained visual perception and fine-grained visual reasoning (Liu et al., 2023b). To investigate both capa-237 bilities of LVLMs, we design a total of 15 tasks in MMDocBench. Specifically, we encompass 238 nine tasks for fine-grained visual perception, including Text Recognition, Table Recognition, Text 239 Localization, Table Cell Localization, Key Information Extraction, Document Forgery Detection, 240 Document Question Answering, Chart Question Answering, and Infographic Question Answering; 241 and for fine-grained visual reasoning, we include six tasks: Arithmetic Reasoning, Logical Reason-242 ing, Spatial Reasoning, Comparison, Sorting and Counting. Further, we include one or multiple 243 sub-tasks for each task to cover more diverse document image types, e.g. research papers, book cov-244 ers, financial reports, scene-text images, receipts, Wikipedia tables, charts, infographics, and other 245 industry documents, leading to a total of 48 sub-tasks in MMDocBench. See a summary in Table 1.

- 246 Step 2: Document Image & QA Pair Preparation. As shown in Table 1, we create BookOCR, 247 Bbox2Text, and Text2Bbox by ourselves and use original task settings for other sub-tasks. In par-248 ticular, we build BookOCR, a text recognition dataset, based on selected document images (i.e., 249 book covers) from OCR-VQA (Mishra et al., 2019). After collecting document images, we use 250 a pre-defined template for text recognition as the question and automatically identify all the OCR 251 content from the image as the ground-truth answer. We build Text2Bbox and Bbox2Text with the 252 same document images, which are selected from DocILE (Šimsa et al., 2023), RVL-CDIP (Harley 253 et al., 2015), DocBank (Li et al., 2020), PubLayNet (Zhong et al., 2019) and PubTabNet (Zhong 254 et al., 2020) to cover a great diversity of document types. The former task requires an LVLM to find the region in the document image given a piece of textual content, while the latter needs the model 255 to identify corresponding text in the document image based on a specified bounding box. For the 256 three sub-tasks, our annotators (6 undergraduate or graduate students majored in computer science) 257 manually create one QA pair for each selected document image. 258
- 259 For other sub-tasks, our annotators manually analyze and select appropriate document images with 260 annotated input-output pairs from the source dataset. After that, the input-output pairs are trans-261 formed into QA pairs following the pre-defined templates. Note that all the document images and QA pairs are selected from the test set of the source dataset except CORD (Park et al., 2019), 262 DUDE (Van Landeghem et al., 2023) and CharXiv (Wang et al., 2024b). CORD has an insufficient 263 number of high-quality document images in its test set, so we select some from its validation set. 264 DUDE and CharXiv have not yet released their test sets, and thus we utilize their validation sets 265 instead. For each sub-task, we at most select 100 document images. 266
- It is worth mentioning that for preparing sub-tasks of fine-grained visual reasoning, we purposely
 select five existing datasets to ensure our MMDocBench includes a great diversity of document
 image types. These datasets include DUDE (Van Landeghem et al., 2023) containing general documents from various industries, WTQ (Pasupat & Liang, 2015) containing table-based documents

309

Main Task	Sub Task	Document Image Type	# Images	# QA Pairs	# Regions
	Fine-	Grained Visual Perception			
Text	TextOCR (Singh et al., 2021)	Scene-Text Images	100	100	100
Recognition	BookOCR (Mishra et al., 2019)	Book Covers	100	100	438
Table	FinTabNet (Zheng et al., 2021)	Financial Reports	100	100	1,864
Recognition	PubTables-1M (Smock et al., 2022)	Scientific Papers	100	100	3,520
Text	Text2Bbox (Simsa et al. (2023) etc.)	Industry Documents	100	100	100
Localization	Bbox2Text (Šimsa et al. (2023) etc.)	Industry Documents	100	100	100
Table Cell	FinTabNet (Zheng et al., 2021)	Financial Reports	100	100	100
Localization	PubTables-1M (Smock et al., 2022)	Scientific Papers	100	100	100
Key	SROIE (Huang et al., 2019)	Receipts	100	303	303
Information	WildReceipt (Sun et al., 2021)	Receipts	100	512	512
Extraction	CORD (Park et al., 2019)	Receipts	100	372	372
Doc Forgery	T-SROIE (Yuxin et al., 2022)	Receipts	100	100	286
Detection	DocTamper (Qu et al., 2023)	Cross-Domain Documents	100	100	129
Document	DocVQA (Mathew et al., 2021)	Industry Documents	100	262	262
OA	WTQ (Pasupat & Liang, 2015)	Wikipedia Tables	100	351	351
Q.1	TAT-DQA (Zhu et al., 2022)	Financial Reports	100	214	214
Chart	ChartQA (Masry et al., 2022)	Cross-Domain Charts	100	104	104
QA	CharXiv (Wang et al., 2024b)	Scientific Charts	100	149	149
Infographic	InfographicVOA (Mathew et al. 2022)	Infographics	100	281	281
QA			100	201	201
	Fine-	Grained Visual Reasoning	12	15	24
	DUDE (Van Landeghem et al., 2023)	General Documents	13	15	34
Arithmetic	WIQ (Pasupat & Liang, 2015)	Wikipedia Tables	54	25	159
Reasoning	AI-DQA (Zhu et al., 2022) CharWise (Warra et al., 2024b)	Financial Table-Text Documents	90	217	433
	Lafa graphic VOA (Mathews et al. 2022)	Scientific Charts	23	23	67
	DUDE (Van Landacham at al. 2022)	Canaral Decomments	10	35	90
	WTO (Decurat & Liong 2015)	Wilingdia Tablaa	10	11	20
Logical	W IQ (Fasupat & Liang, 2013)	Financial Table Tayt Decuments	11	11	41
Reasoning	CharViv (Wang et al., 2022)	Scientific Charts	1	1	12
	InfographicVOA (Mathew et al. 2022)	Infographics	2	2	12
	DUDE (Van Landacham at al. 2022)	Ganaral Documents	29	41	12
Spatial	WTO (Pagupat & Liong 2015)	Wikipadia Tablas	38	41	43
Bassoning	CharViv (Wang et al. 2024b)	Solontific Charts	4	4	12
Reasoning	InfographicVOA (Mathew et al. 2022)	Infographics	17	23	54
	DUDE (Ven Landacham et al. 2022)	Cross Domain Documents	2	2.5	54
	WTO (Pasupat & Liang 2015)	Wikipedia Tables	33	34	74
Comparison	TAT DOA (Zhu et al. 2022)	Financial Table Text Documents	10	10	30
Comparison	CharXiv (Wang et al. 2024b)	Scientific Charts	16	10	50 44
	InfographicVOA (Mathew et al. 2022)	Infographics	13	10	44
	DUDE (Van Landeghern et al. 2023)	General Documents	3	3	
	WTO (Pasunat & Liang 2015)	Wikipedia Tables	6	12	23
Sorting	TAT-DOA (Zhu et al. 2022)	Financial Table-Text Documents	7	12	14
borting	CharXiv (Wang et al. 2024b)	Scientific Charts	15	15	29
	InfographicVOA (Mathew et al. 2022)	Infographics	20	29	57
	DUDE (Van Landeghem et al. 2023)	General Documents	51	55	244
	WTO (Pasupat & Liang, 2015)	Wikipedia Tables	15	15	76
Counting	TAT-DOA (Zhu et al., 2022)	Financial Table-Text Documents	14	14	26
	CharXiv (Wang et al., 2024b)	Scientific Charts	38	40	149
	InfographicVQA (Mathew et al., 2022)	Infographics	44	52	248
		0 1			= .0

Table 1: Taxonomy and statistics of MMDocBench.

from Wikipedia, TAT-DQA (Zhu et al., 2022) containing table-text documents from financial reports, ChartXiv (Wang et al., 2024b) containing chart-based documents from scientific papers, and InfographicVQA (Mathew et al., 2022) containing infographic-based documents.

Step 3: Region Generation. We generate ground truth regions for each QA pair in MMDocBench to
facilitate evaluation. Similar to (Xuan et al., 2024; Chen et al., 2024b), we normalize the coordinates
used to represent the bounding box to the range [0, 1000] w.r.t. the image dimensions.

316 For the tasks regarding fine-grained visual perception, we set the answer's location on the image 317 as the region to be annotated, while for those regarding fine-grained visual reasoning, we annotate 318 the locations of all supporting evidences used to infer the final answer. Specifically, we first obtain 319 the OCR results using Google OCR service for each document image in MMDocBench. For fine-320 grained visual perception tasks, we automatically identify the corresponding value and its bounding 321 box in the OCR result based on the answer. If only one value matches, we use the region of this value as the correct one; otherwise, our annotators manually review and select the appropriate re-322 gion for the answer. For fine-grained visual reasoning tasks, we search for the regions for each 323 supporting evidence if the source dataset already provides the annotation of supporting evidence,



Figure 3: Position distribution of all regions. Figure 4: Area distribution of all regions.

like TAT-DQA (Zhu et al., 2022); for the rest, our annotators manually check supporting evidence and annotate the appropriate region for each one.

Step 4: QA & Region Verification. To ensure high quality of MMDocBench, we further verify the collected data. First, we develop a program to automatically highlight the answer or supporting evidence with its corresponding regions on the document image. Then, different annotators review the rendered document image in three rounds to ensure that the answer and supporting regions to the question are correct.

3.3 STATISTICS AND ANALYSIS

With the above construction pipeline, the resultant MMDocBench contains a total of 2, 400 document images and 4, 338 QA pairs with 11, 353 annotated supporting regions. On average, each question has 10.1 words and around 2.61 supporting regions. Refer to Table 4 in Appendix and Table 1 for more detailed statistics of our MMDocBench.

We analyze the position distribution and area distribution of the annotated supporting regions 356 for each QA pair in our MMDocBench. To compute the position distribution, given a region 357 $[x_1, y_1, x_2, y_2]$, we first calculate the center point by $(\frac{x_1+x_2}{2}, \frac{y_1+y_2}{2})$. Then, we plot all the cen-358 ter points on an image with a dimension of 1000×1000 . As shown in Figure 3, all points are 359 scattered across the entire image, indicating no clear positional bias for supporting regions in MM-360 DocBench. We also analyze the area distribution of all regions to examine their granularity. We first 361 calculate the area of each region and then apply a logarithmic transformation with a base of 10 on 362 the computed area value. As shown in Figure 4, the granularity of regions shows a diversity and the majority of the region areas fall between 1,000 and 10,000, corresponding to regions sized between 10×100 and 100×100 . These analyses well highlight the high quality of our MMDocBench, which 364 is crucial for accurately assessing the capabilities of LVLMs in fine-grained visual understanding. 365

366 367

339 340 341

342

343

344

345

346

347

348 349

- 4 EXPERIMENTS
- 368 369 370

4.1 EXPERIMENTAL SETUP

Evaluated LVLMs. We conduct evaluation experiments with 10 open-source LVLMs and 3 proprietary LVLMs on the proposed MMDocBench. We select those open-source models with strong document understanding capabilities, including LLaVA-V1.6-34B (Liu et al., 2024a), Llava-OV-Chat-72B (Li et al., 2024a), TextMonkey (Liu et al., 2024c), MiniCPM-Llama3-V2.5 (Yao et al., 2024), InternVL2-8B (Chen et al., 2024c), InternVL2-Llama3-76B (Chen et al., 2024c), Qwen2-VL-7B-Instruct (Wang et al., 2024a), mPLUG-DocOwl-1.5-Omni (Hu et al., 2024), mPLUG-Owl3 (Ye et al., 2024), and Ferret (You et al., 2024). For proprietary LVLMs to be evaluated in our experiments, we use Qwen-VL-Max-0809 (Bai et al., 2023) and two versions of the latest OpenAI

Model	Size	Fi Vist	ne-Grain 1al Percep	ed otion	Fi Visu	ne-Graine al Reason	ed 1ing		Overall	
		EM	F1	IOU	EM	F1	IOU	EM	F1	IOU
			Close-S	Source LV	/LMs					
GPT-40	-	62.47	65.18	3.21	70.33	72.56	1.68	66.40	68.87	2.44
GPT-4V	-	52.37	56.95	2.48	50.97	52.22	-	51.67	54.58	1.54
Qwen-VL-Max	-	61.63	65.89	18.60	70.15	72.24	4.27	65.89	69.06	11.44
			Open-S	Source LN	LMs					
InternVL2-Llama3-76B	76B	54.63	58.84	2.30	61.76	64.46	-	58.20	61.65	1.54
LLava-OV-Chat-72b	72B	52.59	57.49	2.28	65.26	67.55	-	58.93	62.52	1.57
mPLUG-DocOwl1.5-Omni	72B	8.18	12.82	1.04	12.24	13.67	-	10.21	13.25	-
LLaVA-V1.6-34B	34B	32.75	39.13	2.66	27.85	31.34	-	30.30	35.23	1.63
Owen2-VL-7B-Instruct	9B	55.15	59.14	15.16	43.07	45.88	2.69	49.11	52.51	8.93
TextMonkey	9B	31.41	35.88	19.22	13.73	14.74	-	22.57	25.31	9.87
MiniCPM-Llama3-V2.5	8B	33.71	40.42	5.93	30.95	33.23	1.45	32.33	36.82	3.69
InternVL2-8B	8B	45.72	51.42	2.01	45.16	48.57	-	45.44	50.00	1.29
mPLUG-Owl3	7B	15.63	19.97	-	10.84	12.65	-	13.24	16.31	-
Ferret	7B	3.96	6.96	4.33	-	-	-	1.98	3.48	2.16

Table 2: The overall performance of different LVLMs on MMDocBench. Best results are marked in
 bold while second-best are underlined. Metric values below 1% are marked with '-'.

GPT (Achiam et al., 2023), i.e. GPT-4o-2024-08-06 and GPT-4-turbo-2024-04-09. These models show a noticeable diversity in respective parameter sizes, visual encoders, and language models.

398 Instruction Design. For each main task in MMDocBench, we manually design an instruction template to guide the LVLM to output the answer and supporting regions in JSON format, e.g., 399 {"answer":"{answer}", "bbox":["{bbox1}","{bbox2}"]}. For fine-grained visual perception tasks, 400 we instruct the LVLM to output the results directly, while for fine-grained visual reasoning tasks, 401 we enable chain-of-thought (CoT) in the instructions. Since some LVLMs cannot follow the in-402 structions well during the test, especially for the region predictions, we revise the instructions based 403 on the setting of these LVLMs to improve their performance. At inference, we apply zero-shot 404 prompting on all the LVLMs to generate responses for each question. After obtaining the response, 405 we develop a program that uses strict regular expressions to extract the predicted answer and the 406 supporting regions for evaluation. 407

Evaluation Metrics. For each question, we use *Exact Match (EM)* and *F1-score* to evaluate the 408 predicted answer, and Intersection over Union (IOU) to assess the predicted region(s). The EM is 409 determined by matching every character of the model's text prediction to the ground truth. If all 410 characters are matched, the *EM* is 1, and otherwise 0. For *F*1-score, we calculate the word-level F1 411 based on the number of words in model prediction, ground truth, and their intersection (Rajpurkar 412 et al., 2018). The Intersection over Union (IOU) is computed between the predicted region and the 413 ground-truth region, taking into account their overlapping area and union area. The scores on the 414 three metrics for each sub-task are computed by taking the mean of corresponding metric scores for 415 all the questions included in this sub-task. To obtain metric scores per task or capability, and overall performance, we calculate the macro average across all corresponding lower-level metric scores. 416

417

419

380 381 382

418 4.2 MAIN RESULTS

We conduct a comprehensive comparison of different LVLMs with our proposed MMDocBench, and 420 show the results in Table 2. We make below key findings. 1) The proposed MMDocBench poses 421 significant challenges to current LVLMs in terms of both answer prediction and region prediction. 422 The best model GPT-40 achieves 66.40% in EM for answer prediction, but only 2.44% in IOU 423 for region prediction. Another close-source model, Qwen-VL-Max, which has comparable answer 424 prediction performance to GPT-40, is the best-performing model for region prediction, with an IOU 425 score of 11.44% only. This highlights substantial challenges of region prediction in MMDocBench. 426 2) There is a large gap in answer prediction performance between open-source and close-source 427 models, while that in their region prediction performance is minor. For answer prediction, GPT-428 40 significantly outperforms the best open-source model Llava-OV-Chat-72b (with 58.93% in EM) 429 by over 12% in EM. For region prediction, open-source models like TextMonkey and Qwen2-VL-7B-Instruct, with IOU scores of 8.93% and 9.87% respectively, perform slightly worse than Qwen-430 VL-Max but outperform GPT-40 by around three times. This could be explained by OpenAI using 431 insufficient samples with visual grounding requirements when training GPT-40 and GPT-4v, leading



Figure 5: Model performance comparison Figure 6: Model performance comparison across all tasks with F1 score. across all tasks with IOU metric.

to a lack of visual grounding capability. Another possible reason is that visual encoders adopted in
GPT-40 and GPT-4v do not effectively support fine-grained visual understanding, which we leave
for future investigation. 3) LVLMs trained on document images with text-grounding requirements,
such as Qwen-VL-Max, Qwen2-VL-7B-Instruct, and TextMonkey, show improvements in region
prediction, while those trained with object-level grounding over natural images, like Ferret, exhibit
no such improvements. This highlights the necessity of establishing a benchmark that supports
visual grounding at various and finer granularities on document images, such as our MMDocBench.

In the following, we further analyze model performance on answer prediction and region prediction regarding fine-grained visual perception and fine-grained visual reasoning, respectively.

Answer Prediction. We make below key findings. 1) GPT-40 consistently beats all other models on
both fine-grained visual perception and fine-grained visual reasoning tasks in terms of EM, demonstrating its superior effectiveness. 2) Larger models, like GPT-40, Qwen-VL-Max, InternVL2Llama3-76B, and Llava-OV-Chat-72b tend to perform better on fine-grained visual reasoning tasks
than fine-grained visual perception tasks, while smaller models, conversely, excel in fine-grained
visual perception tasks over reasoning tasks. This might be because reasoning capabilities improve
significantly as the model size increases.

Region Prediction. We make below key findings. 1) TextMonkey achieves the best results on fine-464 grained visual perception tasks with 19.22% in IOU, but it fails to output supporting regions for fine-465 grained visual reasoning tasks. One possible reason is that TextMonkey only involves perception-466 related samples and instructions for training its grounding ability, resulting in its inability to ground 467 the supporting evidence for reasoning tasks. 2) All LVLMs face significant challenges in predicting 468 supporting regions for fine-grained visual reasoning tasks. The performance of most LVLMs on 469 fine-grained visual reasoning tasks is notably worse than that on fine-grained visual perception tasks. 470 Qwen-VL-Max, which is the best model for region prediction on fine-grained visual reasoning tasks, 471 earns only 4.27% in IOU, indicating the remarkable challenges of this task.

472 473 474

444

445 446 447

4.3 ANALYSIS ON DIFFERENT TASKS

475 We compare the performance of various LVLMs across tasks and present the results in Figure 5 476 and Figure 6. We make below key findings. 1) As shown in Figure 5, GPT-40 and Qwen-VL-477 Max are the best and second-best models across most tasks among all LVLMs, except for Logical 478 Reasoning and Spatial Reasoning, where LLaVA-OV-Chat-72B outperforms all other models. 2) As 479 shown in Figure 6, all LVLMs struggle with the prediction of supporting regions on almost all the 480 tasks, except for Text Recognition, where the best model, Qwen-VL-Max, achieves around 71.2%481 in IOU. Qwen-VL-Max delivers the best results in region prediction across most tasks, except for 482 Forgery Detection, Infographic QA, and Text Localization, where TextMonkey ranks the first. 3) Among all the tasks, Document Forgery Detection, a fine-grained visual perception task, is the 483 most challenging for all LVLMs, with the best result being only around 20.6% in EM for answer 484 prediction and 1.8% in IOU for region prediction; As illustrated in Figure 1, this task requires the 485 model to identify the inconsistent word(s) against other words within the image. In addition, the

Model	Gen	eral	Table-	Based	Table	-Text	Chart	-Based	Infograp	hic-Based
	F1	IOU	F1	IOU	F1	IOU	F1	IOU	F1	IOU
			Close	-Source	LVLMs					
GPT-40	71.21	3.25	71.31	1.39	77.97	2.90	71.53	3.06	74.41	1.38
GPT-4V	58.65	2.43	56.84	-	62.76	2.17	54.89	1.64	64.27	1.14
Qwen-VL-Max	68.45	19.38	67.35	5.23	80.89	9.52	74.18	11.53	79.41	4.39
Open-Source LVLMs										
InternVL2-Llama3-76B	65.72	2.57	60.28	-	68.37	-	68.39	1.38	68.67	-
LLaVA-OV-Chat-72B	61.83	2.65	65.34	-	69.28	1.17	68.39	1.10	69.53	-
mPLUG-DocOwl1.5-Omni	17.06	-	13.35	-	11.84	-	12.99	-	19.69	1.25
LLaVA-V1.6-34B	42.53	2.57	33.42	-	36.65	-	40.62	-	33.96	-
Qwen2-VL-7B-Instruct	55.64	16.57	50.14	2.67	64.12	5.17	55.92	7.57	59.74	2.81
TextMonkey	36.11	15.44	20.17	4.01	24.70	3.34	23.62	5.51	27.35	4.43
MiniCPM-Llama3-V2.5	40.90	6.67	34.97	1.59	41.64	3.26	42.61	-	39.49	3.83
InternVL2-8B	55.24	2.04	52.17	-	53.01	-	57.93	-	52.49	-
mPLUG-Owl3	21.27	-	10.89	-	18.76	-	22.02	-	15.53	-
Ferret	8.61	7.17	3.08	-	4.70	-	5.81	-	5.63	-

Table 3: The performance comparison of LVLMs across different document types. Best results are
 marked in bold; second-best are underlined. Metric values below 1% are marked with a '-'.

answer prediction performance of all LVLMs on Text Localization and Table Cell Localization is poor, underscoring the challenges posed by both tasks. Refer to Section A.2 in Appendix for a more detailed performance comparison of all LVLMs across different tasks.

4.4 ANALYSIS ON DIFFERENT DOCUMENT TYPES

511 Based on the majority of content included, we categorize all the document images in MM512 DocBench into five types: General Document, Table-Based Document, Table-Text Document,
513 Chart-Based Document, and Infographic-Based Document, and ensure document images per sub514 task belong to the same one category. Refer to Table 9 in Appendix for detailed strategies.

We present LVLMs' results across different document types in Table 3, from which we make below findings. 1) For answer prediction, GPT-40 obtains the best performance on General Document and Table-Based Document, while Qwen-VL-Max ranks the first across the other three document types.
2) For region prediction, TextMonkey outperforms all other models on Infographic-Based Document, while Qwen-VL-Max leads on the remaining four document types. 3) In comparison, region prediction on other document types is significantly more challenging than on General Documents, while there is no notable difference in answer prediction across different document types.

5 CONCLUSION

In this work, we introduce the MMDocBench benchmark to comprehensively evaluate LVLMs' fine-grained visual perception and reasoning capabilities via various OCR-free document understanding tasks. In MMDocBench, we carefully design 15 tasks and 48 sub-tasks that require LVLMs to perform deep, fine-grained interpretation of image details to answer each question. To enable a more comprehensive evaluation, we provide annotations of the supporting regions for each question-answer pair to assess whether LVLMs have the abilities to correctly ground their predictions on the associated regions in the image. With MMDocBench, we evaluate various open-source and propri-etary LVLMs, analyzing their performance in fine-grained visual document image understanding. We observe that our MMDocBench presents significant challenges to current LVLMs in both answer and region prediction, with GPT-40 achieving the highest answer prediction score of 66.40% in EM and Qwen-VL-Max achieving the best region prediction score of 11.44% in IOU. Moreover, we find that open-source LVLMs demonstrate competitive performance in region prediction compared to proprietary models, despite a significant gap in answer prediction. We believe MMDocBench can enable a thorough and multi-faceted evaluation of fine-grained visual document understanding of LVLMs, thereby facilitating LVLMs' future advancement.

540 REFERENCES

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
report. *arXiv preprint arXiv:2303.08774*, 2023.

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel
Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan
Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian
Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo
Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language
model for few-shot learning, 2022. URL https://arxiv.org/abs/2204.14198.

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, local ization, text reading, and beyond, 2023.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,
Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin,
Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford,
Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *CoRR*, 2020.

- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi
 Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language
 models? arXiv preprint arXiv:2403.20330, 2024a.
- Wentong Chen, Junbo Cui, Jinyi Hu, Yujia Qin, Junjie Fang, Yue Zhao, Chongyi Wang, Jun Liu, Guirong Chen, Yupeng Huo, et al. Guicourse: From general vision language models to versatile gui agents. *arXiv preprint arXiv:2406.11317*, 2024b.
- 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*, 2024c.
- 571 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong
 572 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning
 573 for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer*574 *Vision and Pattern Recognition*, pp. 24185–24198, 2024d.
- 575
 576
 576
 576
 577
 578
 Chaorui Deng, Qi Wu, Qingyao Wu, Fuyuan Hu, Fan Lyu, and Mingkui Tan. Visual grounding via accumulated attention. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7746–7755, 2018.
- 579 Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu
 580 Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation
 581 benchmark for multimodal large language models, 2024. URL https://arxiv.org/abs/
 582 2306.13394.
- Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang
 Chen, Furong Huang, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. Hallusionbench: An
 advanced diagnostic suite for entangled language hallucination and visual illusion in large vision language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14375–14385, June 2024.
- Adam W Harley, Alex Ufkes, and Konstantinos G Derpanis. Evaluation of deep convolutional nets for document image classification and retrieval. In 2015 13th International Conference on Document Analysis and Recognition (ICDAR), pp. 991–995. IEEE, 2015.
- 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.

611

617

633

634 635

636

637

- Jiaxing Huang and Jingyi Zhang. A survey on evaluation of multimodal large language models.
 arXiv preprint arXiv:2408.15769, 2024.
- Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and C. V.
 Jawahar. Icdar2019 competition on scanned receipt ocr and information extraction. In 2019
 International Conference on Document Analysis and Recognition (ICDAR), pp. 1516–1520, 2019.
- Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document understanding transformer. In *European Conference on Computer Vision*, pp. 498–517. Springer, 2022.
- Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9579–9589, 2024.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei
 Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024a.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image
 pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pp. 19730–19742. PMLR, 2023a.
- Minghao Li, Yiheng Xu, Lei Cui, Shaohan Huang, Furu Wei, Zhoujun Li, and Ming Zhou. Docbank:
 A benchmark dataset for document layout analysis. *arXiv preprint arXiv:2006.01038*, 2020.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 292–305, Singapore, December 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.20. URL https://aclanthology.org/2023.emnlp-main.20.
- 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*,
 2024b.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
 - Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In Advances in Neural Information Processing Systems, pp. 34892–34916. Curran Associates, Inc., 2023a.
 - Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 26296–26306, 2024a.
- 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/.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*, 2023b.
- Yuliang Liu, Zhang Li, Biao Yang, Chunyuan Li, Xucheng Yin, Cheng-lin Liu, Lianwen Jin, and Xiang Bai. On the hidden mystery of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023c.

- Yuliang Liu, Biao Yang, Qiang Liu, Zhang Li, Zhiyin Ma, Shuo Zhang, and Xiang Bai. Textmonkey: An ocr-free large multimodal model for understanding document. *arXiv preprint arXiv:2403.04473*, 2024c.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *International Conference on Learning Representations* (*ICLR*), 2024.
- 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, 2022. Association for Computational Linguistics.
- Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document
 images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*,
 pp. 2200–2209, 2021.
 - Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. Infographicvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1697–1706, 2022.
- Shervin Minaee, Yuri Boykov, Fatih Porikli, Antonio Plaza, Nasser Kehtarnavaz, and Demetri Ter zopoulos. Image segmentation using deep learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(7):3523–3542, 2022.
- Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual question answering by reading text in images. In 2019 international conference on document analysis and recognition (ICDAR), pp. 947–952. IEEE, 2019.
- 674 675 OpenAI. Gpt-4v(ision) system card. 2023.

665

666

667

684

688

689

690

691 692

693

- Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaeheung Surh, Minjoon Seo, and Hwalsuk Lee. Cord: a consolidated receipt dataset for post-ocr parsing. In *Workshop on Document Intelligence at NeurIPS 2019*, 2019.
- Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In Chengqing Zong and Michael Strube (eds.), *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)*, pp. 1470–1480. Association for Computational Linguistics, 2015.
- 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.
 - Chenfan Qu, Chongyu Liu, Yuliang Liu, Xinhong Chen, Dezhi Peng, Fengjun Guo, and Lianwen Jin. Towards robust tampered text detection in document image: New dataset and new solution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5937–5946, 2023.
 - Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, pp. 9, 2019.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don't know: Unanswerable questions for SQuAD. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 784–789. Association for Computational Linguistics, 2018.

702 703 704 705	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. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Pacconition</i> , pp. 13009, 13018, 2024.
706	<i>Vision and Function Recognition</i> , pp. 15009–15010, 2024.
707	Zhongwei Ren, Zhicheng Huang, Yunchao Wei, Yao Zhao, Dongmei Fu, Jiashi Feng, and Xiaojie
709	Jin. Pixellm: Pixel reasoning with large multimodal model. In Proceedings of the IEEE/CVF
700	Conference on Computer Vision and Pattern Recognition, pp. 26374–26383, 2024.
709	v v
710	Stěpán Simsa, Milan Sulc, Michal Uřičář, Yash Patel, Ahmed Hamdi, Matěj Kocián, Matyáš
710	Skalický, Jiří Matas, Antoine Doucet, Mickaël Coustaty, et al. Docile benchmark for document
712	Information localization and extraction. In International Conference on Document Analysis and
713	Recognition, pp. 147–100. Springer, 2023.
714	Amanpreet Singh, Guan Pang, Mandy Toh, Jing Huang, Wojciech Galuba, and Tal Hassner. Textocr:
710	Towards large-scale end-to-end reasoning for arbitrary-shaped scene text. In <i>Proceedings of the</i>
710	IEEE/CVF conference on computer vision and pattern recognition, pp. 8802–8812, 2021.
710	
710	Brandon Smock, Rohith Pesala, and Robin Abraham. Pubtables-1m: Towards comprehensive ta-
719	Die extraction from unstructured documents. In Proceedings of the IEEE/CVF Conference on
720	Computer vision and Pattern Kecognition, pp. 4634–4642, 2022.
721	Hongbin Sun, Zhanghui Kuang, Xiaoyu Yue, Chenhao Lin, and Wayne Zhang, Spatial dual-modality
722	graph reasoning for key information extraction. arXiv preprint arXiv:2103.14470, 2021.
723	
724	Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide
725	shut? exploring the visual shortcomings of multimodal llms. In Proceedings of the IEEE/CVF
726	<i>Conference on Computer Vision and Pattern Recognition</i> , pp. 9568–9578, 2024.
727	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
728	Lacroix Baptiste Rozière Naman Goval Fric Hambro Faisal Azhar et al Llama: Open and
729	efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023.
730	
731	Jordy Van Landeghem, Rubèn Tito, Łukasz Borchmann, Michał Pietruszka, Pawel Joziak, Rafal
732	Powalski, Dawid Jurkiewicz, Mickaël Coustaty, Bertrand Anckaert, Ernest Valveny, et al. Docu-
733	ment understanding dataset and evaluation (dude). In <i>Proceedings of the IEEE/CVF International</i>
734	Conference on Computer Vision, pp. 19528–19540, 2023.
735	Guangzhi Wang, Yixiao Ge, Xiaohan Ding, Mohan Kankanhalli, and Ying Shan. What makes for
736	good visual tokenizers for large language models? <i>arXiv preprint arXiv:2305.12223.</i> 2023a.
737	
738	Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou,
739	Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a
740	simple sequence-to-sequence learning framework. In International conference on machine learn-
741	<i>ing</i> , pp. 23318–23340. PMLR, 2022.
742	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Kegin Chen, Xuejing Liu
743	Jialin Wang, Wenbin Ge, et al. Owen2-vl: Enhancing vision-language model's perception of the
744	world at any resolution. arXiv preprint arXiv:2409.12191, 2024a.
745	
746	Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,
747	Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang.
748	Cogvinit: visual expert for pretrained language models, 2023b.
749	Zirui Wang, Mengzhou Xia, Luxi He. Howard Chen. Yitao Liu. Richard Zhu. Kaiou Liang.
750	Xindi Wu, Haotian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev Arora, and Dangi Chen.
751	Charxiv: Charting gaps in realistic chart understanding in multimodal llms. arXiv preprint
752	arXiv:2406.18521, 2024b.
753	
754	Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. Lulm abubi. A comprehensive qualitation handwards for large
755	vision-language models. arXiv preprint arXiv:2306.09265, 2023.

784

792

796

797

798

- Shiyu Xuan, Qingpei Guo, Ming Yang, and Shiliang Zhang. Pink: Unveiling the power of referential comprehension for multi-modal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13838–13848, 2024.
- 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, 2024. URL https://arxiv.org/
 abs/2408.01800.
- Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Yuhao Dan, Chenlin Zhao, Guohai Xu, Chenliang Li, Junfeng Tian, et al. mplug-docowl: Modularized multimodal large language model for document understanding. *arXiv preprint arXiv:2307.02499*, 2023a.
- Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Guohai Xu, Chenliang Li, Junfeng Tian, Qi Qian, Ji Zhang, Qin Jin, Liang He, Xin Lin, and Fei Huang. UReader: Universal OCR-free visually-situated language understanding with multimodal large language model. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 2841–2858, Singapore, 2023b. Association for Computational Linguistics.
- 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 multi-modal large
 language models. *arXiv preprint arXiv:2408.04840*, 2024.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023c.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on
 multimodal large language models, 2024. URL https://arxiv.org/abs/2306.13549.
- Kaining Ying, Fanqing Meng, Jin Wang, Zhiqian Li, Han Lin, Yue Yang, Hao Zhang, Wenbo Zhang,
 Yuqi Lin, Shuo Liu, et al. Mmt-bench: A comprehensive multimodal benchmark for evaluating
 large vision-language models towards multitask agi. *arXiv preprint arXiv:2404.16006*, 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. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=2msbbX3ydD.
- Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, pp. 69–85. Springer, 2016.
 - Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. arXiv preprint arXiv:2308.02490, 2023.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9556–9567, 2024a.
- Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Ming Yin, Botao Yu, Ge Zhang, et al. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. *arXiv preprint arXiv:2409.02813*, 2024b.
- WANG Yuxin, ZHANG Boqiang, XIE Hongtao, and ZHANG Yongdong. Tampered text detection via rgb and frequency relationship modeling. *Chinese Journal of Network and Information Security*, pp. 29, 2022.

810	Yanzhe Zhang, Ruivi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Divi Yang, and Tong Sun.	
811	Llavar: Enhanced visual instruction tuning for text-rich image understanding. arXiv preprint	1
812	arXiv:2306.17107, 2023.	

- 813
 814 Xinyi Zheng, Doug Burdick, Lucian Popa, Peter Zhong, and Nancy Xin Ru Wang. Global table
 815 extractor (gte): A framework for joint table identification and cell structure recognition using
 816 visual context. *Winter Conference for Applications in Computer Vision (WACV)*, 2021.
- Xu Zhong, Jianbin Tang, and Antonio Jimeno Yepes. Publaynet: largest dataset ever for document layout analysis. In *2019 International Conference on Document Analysis and Recognition* (*ICDAR*), pp. 1015–1022. IEEE, 2019.
- Xu Zhong, Elaheh ShafieiBavani, and Antonio Jimeno Yepes. Image-based table recognition: data, model, and evaluation. In *European conference on computer vision*, pp. 564–580. Springer, 2020.
- Beyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models, 2023. URL https://arxiv.org/abs/2304.10592.
 - Fengbin Zhu, Wenqiang Lei, Fuli Feng, Chao Wang, Haozhou Zhang, and Tat-Seng Chua. Towards complex document understanding by discrete reasoning. In *Proceedings of the 30th ACM International Conference on Multimedia*, pp. 4857–4866, 2022.
- Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7w: Grounded question answering in images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4995–5004, 2016.

864 A APPENDIX 865

868 869

870 871

866 A.1 KEY STATISTICS 867

We present some key statistic about MMDocBench in Table 4.

Table 4: Key statistics of MMDocBench.

872	Statistic	Number
873	Total Number of Tasks	15
874	- # Visual Perception Tasks	9
875	- # Visual Reasoning Tasks	6
876	Total Number of Sub Tasks	48
878	- # Visual Perception Sub-Tasks	10
879	- # Visual Reasoning Sub-Tasks	29
880	Number of Existing Datasets Involved	2)
881	Total Number of Images	$\frac{21}{2400}$
882	Total Number of OA Dains	2,400
883	Total Number of QA Pairs	4,338
884	Total Number of Regions	11,353
885	Avg. Number of Words per Question	16.14
886	Avg. Number of Words per Answer	4.08
887	Avg. Number of Regions per Question	2.61
888		
800		
891		
892		
893		
894		
895		
896		
897		
898		
899		
900		
901		
902		
903		
904		
905		
907		
908		
909		
910		
911		
912		
913		
914		
915		
916		
917		

A.2 DETAILED PERFORMANCE OF LVLMS ACROSS DIFFERENT TASKS

Table 5: Model performance comparison across different fine-grained visual perception tasks using F1 score. TXR:Text Recognition, TBR:Table Recognition, TXL:Text Localization, TCL:Table Cell Localization, KIE:Key Information Extraction, DFD:Document Forgery Detection, DQA:Document Question Answering, CQA:Chart Question Answering, IQA:Infographic Question Answering.

Model	Fine-Grained Visual Perception								
	TXR	TBR	TXL	TCL	KIE	DFD	DQA	CQA	IQA
GPT-40	95.66	73.53	10.81	42.14	89.23	20.44	88.76	85.18	80.90
GPT-4V	88.28	73.33	14.12	30.04	82.45	1.99	79.08	72.31	70.93
Qwen-VL-Max	91.85	70.05	29.09	39.59	90.66	6.74	90.23	88.77	86.00
InternVL2-Llama3-76B	90.82	75.18	12.90	33.25	85.97	9.24	74.06	80.53	67.64
LLaVA-OV-Chat-72B	90.58	68.12	5.74	26.55	79.67	4.49	81.09	84.16	76.99
mPLUG-DocOwl1.5-Omni	33.90	-	0.80	2.19	7.37	2.25	15.03	16.73	24.32
LLaVA-V1.6-34B	77.68		4.43	11.97	67.14	2.92	54.88	55.27	38.72
Qwen2-VL-7B-Instruct	89.04	76.11	24.57	28.00	81.63	0.67	79.23	78.21	74.76
TextMonkey	88.13	-	26.10	8.70	59.98	5.50	53.12	42.66	38.71
MiniCPM-Llama3-V2.5	78.40	47.40	11.91	18.27	55.52	1.01	51.40	56.01	43.83
InternVL2-8B	79.09	74.93	6.51	21.70	78.53	5.76	67.21	72.26	56.81
mPLUG-Ow13	45.73	5.31	0.12	7.67	38.20	1.59	28.07	36.01	17.04
Ferret	29.19	-	5.32	0.57	0.69	0.93	7.56	5.81	5.63

Table 6: Model performance comparison across different fine-grained visual perception tasks using IOU. TXR:Text Recognition, TBR:Table Recognition, TXL:Text Localization, TCL:Table Cell Localization, KIE:Key Information Extraction, DFD:Document Forgery Detection, DQA:Document Question Answering, CQA:Chart Question Answering, IQA:Infographic Question Answering.

Model			Fir	ne-Grain	ed Visual	Percepti	on		
	TXR	TBR	TXL	TCL	KIE	DFD	DQA	CQA	IQA
GPT-40	14.64	0.93	2.09	2.47	2.04	0.25	2.06	3.20	1.18
GPT-4V	8.66	0.73	1.55	1.30	4.42	0.02	1.66	2.40	1.57
Qwen-VL-Max	70.80	5.59	10.91	6.40	33.22	1.12	16.79	16.96	5.60
InternVL2-Llama3-76B	14.48	0.61	1.20	0.50	1.45	0.30	0.62	0.90	0.62
LLaVA-OV-Chat-72B	13.48	0.25	1.18	0.53	2.60	0.05	1.03	0.72	0.65
mPLUG-DocOwl1.5-Omni	3.49	-	1.01	0.03	0.30	0.12	0.93	0.14	2.27
LLaVA-V1.6-34B	11.01		0.55	1.49	5.45	0.32	0.94	0.35	1.13
Qwen2-VL-7B-Instruct	65.55	2.23	7.08	3.12	30.63	0.00	12.27	11.28	4.28
TextMonkey	67.54	-	43.64	2.39	24.30	1.80	14.72	10.48	8.15
MiniCPM-Llama3-V2.5	12.94	1.09	9.18	0.68	13.03	0.03	9.26	0.40	6.74
InternVL2-8B	11.86	0.42	1.03	1.00	2.34	0.07	0.58	0.48	0.34
mPLUG-Owl3	2.03	0.21	0.02	0.31	0.39	0.01	0.15	0.12	0.1
Ferret	29.12	-	1.47	0.98	0.98	0.29	0.83	0.74	0.2

974							
975	Model		Fine-Gra	ined Visual Reasonir	ıg		
976		Arithmetic Reasoning	Logical Reasoning	Spatial Reasoning	Comparison	Sorting	Counting
977	GPT-40	76.52	73.55	67.50	74.64	73.86	69.28
011	GPT-4V	53.36	49.25	54.08	49.18	66.57	40.89
978	Qwen-VL-Max	82.37	70.27	66.80	69.84	76.70	67.45
979	InternVL2-Llama3-76B	65.77	70.28	65.44	62.27	64.83	58.17
010	LLaVA-OV-Chat-72B	65.91	75.09	77.18	71.51	62.88	52.74
980	mPLUG-DocOwl1.5-Omni	4.92	11.17	22.14	18.76	17.30	7.72
981	LLaVA-V1.6-34B	28.37	44.61	43.46	37.12	18.10	16.36
001	Qwen2-VL-7B-Instruct	50.42	43.61	53.42	48.93	38.26	40.62
982	TextMonkey	2.11	10.91	20.38	25.68	14.37	14.99
983	MiniCPM-Llama3-V2.5	26.73	32.37	44.40	44.54	39.57	11.78
000	InternVL2-8B	48.19	56.64	60.08	49.73	41.30	35.50
984	mPLUG-Owl3	2.81	25.34	16.79	22.23	5.39	3.36
985	Ferret	0.00	0.00	0.00	0.00	0.00	0.00

Table 7: Model performance comparison across different fine-grained visual reasoning tasks usingF1.

Table 8: Model performance comparison across different fine-grained visual reasoning tasks using IOU.

	Model		Fine-Gra	ained Visual Reasonir	ıg		
1		Arithmetic Reasoning	Logical Reasoning	Spatial Reasoning	Comparison	Sorting	Counting
2	GPT-40	0.66	1.70	2.81	0.70	1.85	2.37
	GPT-4V	0.41	0.44	0.98	0.44	0.42	0.91
3	Qwen-VL-Max	3.14	4.68	7.15	4.90	2.63	3.14
4	InternVL2-Llama3-76B	0.50	1.73	1.43	0.18	0.04	0.86
-	LLaVA-OV-Chat-72B	0.42	0.91	1.92	0.31	0.58	1.08
5	mPLUG-DocOwl1.5-Omni	0.22	0.02	0.21	0.13	0.06	0.17
6	LLaVA-V1.6-34B	0.17	1.39	1.02	0.28	0.43	0.32
-	Qwen2-VL-7B-Instruct	1.74	4.66	4.99	1.92	1.18	1.66
1	TextMonkey	0.00	0.00	1.12	0.52	1.41	0.00
8	MiniCPM-Llama3-V2.5	1.16	1.49	2.37	2.64	0.28	0.77
0	InternVL2-8B	0.11	1.06	1.43	0.21	0.18	0.45
9	mPLUG-Owl3	0.00	0.04	0.71	0.00	0.02	0.00
00	Ferret	0.00	0.00	0.00	0.00	0.00	0.00

1026 A.3 DOCUMENT CATEGORY

We follow the strategy in Table 9 to divide the document images in each sub-task into five categories,
i.e., General Document, Table-Based Document, Table-Text Document, Chart-Based Document and
Infographic-Based Document.

Table 9: The document image category for each sub task

Task	Sub Task	Category
Text Recognition	TextOCR	General Document
Text Recognition	OCR-VQA	General Document
Table Descention	FinTabNet	Table-Based Document
Table Recognition	PubTables-1M	Table-Based Document
Tout Localization	Text2Bbox	General Document
lext Localization	Bbox2Text	General Document
	FinTabNet	Table-Based Document
Table Cell Localization	PubTables-1M	Table-Based Document
	SROIE	General Document
Key Information Extraction	WildReceipt	General Document
.,	CORD	General Document
	T-SROIE	General Document
Document Forgery Detection	DocTamper	General Document
	DocVOA	General Document
Document Question Answering	WTO	Table-Based Document
Document Question 7 inswering	TAT-DOA	Table-Text Document
	ChartOA	Chart Based Document
Chart Question Answering	CharViv	Chart Based Document
Infographic Question Answering	InfographicVOA	Inforgraphic Based Documen
Intographic Question Answering		Canaral Dacument
	DUDE	Table Deced Decument
A with we off a December a		Table-Based Document
Arithmetic Reasoning	IAI-DQA CharVis	Chart Desed Desument
	UnarAiv Information VOA	Charl-Based Document
	Intographic VQA	Inforgraphic-Based Documen
	DUDE	General Document
I ' 1D '	WIQ	Table-Based Document
Logical Reasoning	IAI-DQA	Table-Text Document
	CharXiv	Chart-Based Document
	InfographicVQA	Inforgraphic-Based Documen
	DUDE	General Document
Spatial Reasoning	WIQ	Table-Based Document
Spana reasoning	CharXiv	Chart-Based Document
	InfographicVQA	Inforgraphic-Based Documen
	DUDE	General Document
	WTQ	Table-Based Document
Comparison	TAT-DQA	Table-Text Document
	CharXiv	Chart-Based Document
	InfographicVQA	Inforgraphic-Based Documen
	DUDE	General Document
	WTQ	Table-Based Document
Sorting	TAT-DQA	Table-Text Document
-	CharXiv	Chart-Based Document
	InfographicVQA	Inforgraphic-Based Documen
	DUDE	General Document
	WTO	Table-Based Document
Counting	TAT-DOA	Table-Text Document
- 0	CharXiv	Chart-Based Document
	InfographicVOA	Inforgraphic-Based Documen
	mographic VQA	morgraphic-based Documen

1080 A.4 SAMPLES FOR EACH TASK IN MMDOCBENCH

1082 A.4.1 TEXT RECOGNITION 1083



- Question: Could you identify and read the text present in the provided image?
- Answer: Chaos Computer Club
- Bounding Box: [339, 362, 856, 498]

A.4.2 TEXT LOCALIZATION

1102	
1103	BIOGRAPHICAL SKETCH OF CO-INVESTIGATOR :
1104	Name : Uday Gadgil,M.D.
1105	Place of Birth : Bombay, India
1106	Work Address : Dept of Cardiology
1107	Duarte, CA 91010.
1108	Work Phone : (818) 359 - 8111 x 2491
1109	Education :
1110	University Dates Major Degree
1111	University of Bombay '66 - '70 Medicine M.D.,
1112	MEMBERSHIP :
1113	Member of Los Angeles Society of Echocardiography
1114	American Heart Association , Greater Los Angeles ,Affiliate Fellow , American College of cardiolgy, 1978 to present
1115	POSITIONS HELD :
1116	7/83 - present : Staff Cardiologist City of Hope medical Center
1117	7/82 = 6/82
1118	Mount Sinai Medical Center
1119	RESEARCH PUBLICATIONS : About 7 publications
1120	indu - pullidu - nodu - pullidulons .
1121	
1122	• Text2Bhox Question: Can you find where "Miami Florida" appears in the image and
1123	send back its position?
1124	• Downding Dow $[A06, 644, 576, 659]$
1125	• Domaing Dox . [400, 044, 570, 058]
1126	• Bbox2Text Question : For the provided image and bounding box [406, 644, 576, 658], can
1127	you extract and return the text contained within the specified area?
1128	• Text: Miami, Florida.
1129	
1130	
1131	
1132	
1133	

1134 A.4.3 TABLE RECOGNITION 1135 1136 1137 1138 TABLE 6: MSC stemness genes. pr 1139 N. 1140 Total Promoter Inside Downstream ZI 1141 Total genes CGIs CGIs CGIs CGIs 1142 K_{i} 1143 42 49 17 31 1 by 1144 Unvaried 29 1 10 18 da 1145 th Unmet wave 6 0 16 10 1146 1147 R'. Met wave 4 1 3 0 1148 pł 1149 ac 1150 TABLE 7: MSC differentiation genes. El1151 1152 1153 • Question: Could you break down the table into its cell components? 1154 [[Total genes, Total\nCGIs, Promoter\nCGIs, Inside\nCGIs, Down-• Answers: 1155 stream\nCGIs], [42, 49, 17, 31, 1], [Unvaried, 29, 10, 18, 1], [Unmet wave, 16, 6, 10, 1156 0], [Met wave, 4, 1, 3, 0]] 1157 • Bounding Box: [[[100,252,236,347], [320,217,385,382], [445,217,559,382], 1158 [614,217,687,382], [739,217,896,382]], [[100,394,130,482], [339,394,369,482], 1159 [[100,494,211,582], [489,394,513,482], [638, 394, 663, 482],[812,394,823,482]], 1160 [339,494,366,582], [489,494,516,582], [638,494,663,582], [812,494,823,582]], 1161 [[100,588,252,676], [339,588,366,676], [494,588,510,676], [638, 588, 663, 676],1162 [809,588,823,676]], [[100,688,214,776], [345,688,361,776], [497,688,508,776], [644,688,657,776], [809,688,823,776]]] 1163 1164 A.4.4 TABLE CELL LOCALIZATION 1165 1166 0.00 E. 1 1167 1168 1169 1170 Supplem. n [embryos (dams)] Normal Malformed 1171 a. None 155 (22) 0% 100% 1172 b.bC^H 144 (15) 2% 98% 1173 c. apoAL³⁴¹ 37.(6) 0% 100% 1174 1175 Table 3. Phenotype distribution of embryos at 14.5 dpc from B 1176 vitamin A-deficient diet during pregnancy and under various re 1177 1178 1179 • **Question**: What is stored in the cell at the intersection of row 1 and column 3? 1180 • Answer: Normal 1181 1182 • Bounding Box: [594, 283, 684, 373] 1183 1184 1185 1186 1187

1188 A.4.5 KEY INFORMATION EXTRACTION

1190 1191 1192 1193 1194 1195 1196 1197 1198 1199 1200 1201 1202 1203 1204 1205 1206 1207 1208 1209 1210 1211 1212 1213 A.4.6 1217 1218 1219 1220 1221	Question: Please perform detection and recognition of the text that has been categorized under the "total amount" label. P. nswer: 50.002 Bounding Box: [685, 740, 799, 762] DOCUMENT FORGERY DETECTION
1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238 1239	• over 30 frames • the camera im- • inherent scale \$ lobal 3D results \$ d not supported \$ n metrically ac- etric length of a $Temporal Smoot ations from constan E_{temp}(t)where the gradientsmined using finite (the territory of the set of the territory of the territory of the set of the territory of the territory of the territory of the set of the territory of territo$
1240 1241	 Answer: results Bounding Box: [220, 707, 380, 755]





1296 A.4.9 INFOGRAPHIC QUESTION ANSWERING

1298			
1299			
1300			
1301			
1302			
1303			IVIENTAL TEALTTCKISES If you suspect that someone you know is considering suicide or otherwise harming
1304			themselves, seek profesional help as soon as possible and don't leave them alone. CALL 911 IN CASE OF IMMEDIATE EMERGENCY
1305			COMMON SIGNS OF A MENTAL HEALTH CRISIS
1306			
1307			Difference in Initiability Isolation Hopelessness Neglect of demonstration in the second
1308			a lowed on starts the conversion survey from a starts. If you have back and the starts and the s
1309			time and place maybe available to avail this person Be a good listener V Encourage them to make the cell V Avaid distructions
1310			MENTAL HEALTH CRISIS HOTLINES
1311			For drug, alcohel and For these of Native descent: For subidal individuals: gambing addicts: Fors: Nations and holt Rope: Centre for Subida Connectionary for Wellows (Reg Loss Prevention
1312			For claim For cating disorder sufference For claim Streams For cating disorder sufference For claim suff
1313			1-302-668-8686 1-366-638-4820 1-355-450-4506
1314			Ontario.cmha.ca/documents/are-you-in-crisis/
1315			For more information about mental health first aid, read our monthly newsletter at aja ca www.liptomerg.sexeed.intenderg.httl:/ajacates.html
1316			
1317	•	Question: What is th	ne contact number for children and teens suffering from mental crisis,
1318		1-833-456-4566, 1-8	00-668-6868, or 1-866-633-4220?
1319	•	Answer: 1-800-668-	6868
1320		Bounding Box: [147	820 331 8451
1321		Dounding Dox. [147	, 027, 551, 045]
1322	A 4 10	COUNTING	
1323	11.1.10	Cooninto	
1324			
1325			Employer Stock Parchase Plan Highle employees are effected shared from the Steining period, which consists of four consecutive 6-meeth parchase
1326			periods. Europhyses may produce an limited mather of shares of the Company's stock at a discount of up 15% of the bisisse of the mather stude at the high theorymin of the efformation of the output of theorymin of the tentions of the stars was matched to increase the high storing stars of the company stock. As of April 35, 2019. "Timilion dates ware would be formance. The biological period numerics arisely storing and and the ESPP (m
1327			millons): Vier Talad Street Auto: April X, 2017 Street Auto: April X, 2017
1328			Interest function on a last set of the set o
1329			Stock-based compensation expense is included in the consolidated statements of operations as follows (in million): Year Faird April 1984 Stock 1984 Stock 1984
1330			Coard forwhold treatments 5 4 5 3 5 4 Coard forwhold treatmentsmanned other services revenues 10 10 13 Sales and machining 67 60 84 Remember and down 67 60 60
1331			Concernance 79 31 35 Total ance-based compensation expranse \$ 158 \$ 161 \$ 195 Income the based compensation \$ 155 \$ 161 \$ 195
1332			As of April 26, 2019, busi anrecognized compensation expense related to our equity search was 5215 million, which is expected to be recognized as a sample-in-basis over a weighted-average remaining service prioris of 2.1 years. Faduated - commonion
1333			 The valuation of RSUs and ESPP purchase rights and the underlying weighted-average assumptions are summatized as follows: Yur Todat
1334			April 23.00 April 23.00 April 23.00 Risk-levi intrest nize 2.0% 1.4% 1.0% Expected divided joid 2.0% 2.0% 3.1%
1335			Weighted-wronge fuir value par have ganzed \$ 63.40 \$ 39.74 \$ 24.99 ESPP: Expended turnin years 1.2 1.2 1.2
1336			Risk-free interest rate 2.4% 1.4% 0.8% Expected dowling bit 31% 2.2% 3.0% Expected dowling bit 2.4% 2.0% 3.1% Weinhold-more first part of the remaind 5 10.4 7.85
1337			Streed, Reparations Program
1338			stek regelstes program. Under des program, which we may support or disconting at any time, see energy performable adhers of our continuing control des programs. The energy control of the output of the performance of the pe
1339			76
1340			
1341			
1342	•	Question: How man	y years did Shares issued under the ESPP exceed \$2 million?
1243	•	Reasoning Type: co	unting
1344	-	Answer: 3	
1346	•	Anower, J	
1347	•	Bounding Box: []	
1348	•	Supporting Evidence	:e:
1349		– Text : [4, 3, 4]	
		- Bounding Box:	[[795, 225, 807, 236], [673, 225, 685, 236], [918, 225, 929, 236]]

1350 A.4.11 ARITHMETIC REASONING 1351

1352			Churches in Le	evanger		
1353		Parish	Church name	Location	Year built	
1354		Alstadhaug	Alstadhaug Church	Alstadhaug	1180	-
1355		Ekne	Ekne Church	Ekne	1893	-
1356		Levanger	Bamberg Church	Levanger	1902	-
1357		Markabygd	Markabygda Church	Markabygd	1887	
1358		Okkenhaug	Okkenhaug Chapel	Okkenhaug	1893	
1359		Ytterøy	Ytterøy Church	Ytterøya	1890	
1360		Åsen	Åsen Church	Åsen	1904]
1361						
1001			0.1.5			
1302	• Question: how m	any years	s after the levan	ger churc	h was bu	ilt was the bamberg ch
1363	built?					
1364	Reasoning Type:	arithmetic	c			
1365	• Answer [.] 96					
1366						
1367	• Bounding Box: []					
368	 Supporting Evide 	ence:				
369	– Text · [1998	19021				
370	_ Rounding D	w. [[2/2]	506 018 5511	[8/3 /07	017 15	211
371	- Bounding Bo	JA. [[043,	500, 910, 551],	1043,407	, 717, 43	<u>ل</u> [ک
372	A 4 12 Logicit Drive	anna				
373	A.4.12 LOGICAL REASO	JNING				
374						
375			laborator	Tootio	a Drac	0.00
376		רשוער	Laborator	yrestin	y Proc	622
377	Below is a simplified 10-step proc This is a very	ess to walk you th precise process v	rough how COVID-19 testin with no room for error, so it 1	ng is done by mee nust performed b	lical laboratory t y skilled individ	echnologists (referred to as MLTs). uals like MLTs.
378						
370		D	uring extraction, MLTs make co	mplex		
380	*		alculations for a "master mix" t mplify viral RNA.			Results are transcribed into a report
204	The second se	A				ordering healthcare provider and local health unit immediately
201	Sample is taken			10		
202	from a patient's Nasopharyngeal		UUU	extraction,	-	
303	or throat, using a swab.		*	carefully pipette small	44	
304				amounts of the sample into		8
385		× sr	ecimen reaches	the master mix	- {♥ }∧f	iter completion,
386	A AND		e iab. Specially ained MLTs extract	C C	m ct	eticulously necked by
387	Swab goes into	ex	straction. This takes		M	LTs to ensure ey pass Quality Doctor or
388	a tube and is prepared for the				Cim	easures for test
389	laboratory.	-7	U		pe	imple controls.
390			Once mi put throu	xed, samples are ugh an analyzer in a	1 chain	appropriate treatment.
391			reaction reaction	or PCR. This detects	ts if a	
392	3 Sample	e is sent to the lab, onsite or at a testin	g NOTE: step	os 4-7 take a total of :	5 hours	CSMLS CSLM Canadian Society for Wadual Laboratory Science
393	facility	. Transit times vary.				Bocress considienne de science de laborateire médiçai
394						
395	· Onestion William	aton dana	too the arrality		out of 1	toot nonformer of form OO
396	• Question: which	step deno	tes the quality a	ssurance p	part of the	test performed for CO
397	19?					
398	Reasoning Type:	logical				
300	• Answer: 8					
100	Describer D	106 506 7	44 5691			
1404	• Bounding Box: [7	20,536,74	44,308]			
401	 Supporting Evide 	ence:				
1402	– Text : [they n	ass Qualit	v Control (OC)	1		
1403	– Bounding Bo	ox: [[712,	685,811,737]]			

		AL C		2.		
		ияли		Durante		
	Casar Free Marr	чибе им. Г. Ибратимова АН РТ Федала	Википедия Ирекле энциклопедия	PRING	Veparese WIKIMEE	OSTERREICH FOUNDATION
			DIPLO	МА		
			This is to confirm	a that		
		REPLACE	this field with YOUR na	me or account	and PRINT	
	made	the course of the	contribution into developing eighth annual CEE Sprin	g Tatar language o g @ tt.wikipe	content on the Int dia.org marath	ternet 10n
		taking place o	n Tatar Wikipedia from Mar	ch 21st through M	ay 31 st , 2022.	
	Co-founder and Executive Director of Tatar Internet Development Foundation	Member of Tatarstan Academy of Science Scientific Director of Applied	Corresponding member of es, Tatarstan Academy of Science Semiotics Director of G.Ibragimov Insti	Tatar Wikipedia / 2018 Wikimediar tute of Member of Tatar	Administrator 1 of the Year & Turkic User Groups	Director of «Wikimedia RU» Non-Profit Partnership for the Advancement of Encyclopedic
	(Jewels of Knowledge, Literary Marathon et al. projects)	Institute, Doctor of Technical	Sciences Language, Literature and Arts Doctor of Philology	 Local Contest org representative in 	zanizer, ttWP Tatar Portals Union	Knowledge
	(Tatar online content) Member of Tatar Portals Union	Conference Co-Founder and Member, Selet Foundation E General Wiki-Collaboration	Board Presidential Commission for 1 Director Preservation and Strengthenin initiator Tatar language use	the Presidential Com ag of Preservation and language use	mission for the Strengthening of Tatar	Wikimedia Languages of Russia volunteers community
	Rail M. Gataullin	Professor Djavdet Sh. Suley	manov Professor Kim M. Minnullin	Farhad N. Fatkul	lin	Vladimir V. Medeyko
		ON			1	
	Tansonador		200 Afring		from growthe	
	8 Finalists, their qualifyin	ng contribution and aware	ds are public at https://w.wiki/5Heq			
•	Ouestion . Where	is the larg	est font word p	rinted at th	he top or	bottom of the
	Descening Type:	spotiol			it top of	
•	Keasoning Type.	spatial				
•	Answer: Top					
	1					
•	Bounding Box: []					
•	Bounding Box: [] Supporting Evide	ence:				
•	Bounding Box: [] Supporting Evide – Text: [DIPL0	ence: DMA1				
•	Bounding Box: [] Supporting Evide – Text: [DIPLC – Bounding Bo	ence: DMA] ox: [[382.	222. 610.268]]			
•	Bounding Box: [] Supporting Evide – Text: [DIPLC – Bounding Bo	ence: DMA] ox: [[382,	222, 610,268]]			
• • A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] ox: [[382,	222, 610,268]]			
• • A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] ox: [[382,	222, 610,268]]			
• • A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLC – Bounding Bo COMPARISON	ence: DMA] DX: [[382,	222, 610,268]] Churches in Le	vanger	Yana kuilt	1
• • A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] ox: [[382, Parish Alstadhaug	222, 610,268]] Churches in Le Church name Alstadhaug Church	Location	Year built 1180	
• • A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] ox: [[382, Parish Alstadhaug Ekne	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church	Location Alstadhaug Ekne	Year built 1180 1893	
• • A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] ox: [[382, Maisadhaug Ekne Levanger	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church	Location Alstadhaug Ekne Levanger	Year built 1180 1893 1902	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] DX: [[382, Parish Alstadhaug Ekne Levanger Markabugd	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Bamberg Church	Levanger Location Alstadhaug Ekne Levanger Markabud	Year built 1180 1893 1902 1998	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] DX: [[382, Parish Alstadhaug Ekne Levanger Markabygd Okkenhaug	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Bamberg Church Markabygda Church Okkenhaug Chapel	Location Alstadhaug Ekne Levanger Levanger Markabygd Okkenhaug	Year built 1180 1893 1992 1998 1887 1893	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] DX: [[382, Markabaug Ekne Levanger Markabygd Okkenhaug Ytterøy	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Bamberg Church Markabygda Church Okkenhaug Chapel Ytterøy Church	Levanger Levanger Levanger Levanger Markabygd Okkenhaug Ytterøya	Year built 1180 1893 1902 1998 1887 1893 1890	
• • A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] DX: [[382, Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Bamberg Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church	Location Alstadhaug Ekne Levanger Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1902 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] ox: [[382, Marish Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Bamberg Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church	Location Alstadhaug Ekne Levanger Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1902 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] ox: [[382, Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Bamberg Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church	Vanger Location Alstadhaug Ekne Levanger Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1902 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON	ence: DMA] DX: [[382, Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Bamberg Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church 1930 above/bele	Location Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1902 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON Question: was the Reasoning Type:	ence: DMA] DX: [[382, Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church 1930 above/belo	Levanger Levanger Levanger Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1992 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON Question: was the Reasoning Type: Answer: above	ence: DMA] ox: [[382, Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Bamberg Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church 1930 above/belo	Location Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1902 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON Question: was the Reasoning Type: Answer: above Bounding Box: []	Parish Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Bamberg Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church 1930 above/belo	Location Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1902 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON Question: was the Reasoning Type: Answer: above Bounding Box: []	ence: DMA] ox: [[382, Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Bamberg Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church 1930 above/belo	Location Alstadhaug Ekne Levanger Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1902 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON Question: was the Reasoning Type: Answer: above Bounding Box: [] Supporting Evide	ence: DMA] DX: [[382, Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen e finish in compariso	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church 1930 above/belo	Location Alstadhaug Ekne Levanger Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1902 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON Question: was the Reasoning Type: Answer: above Bounding Box: [] Supporting Evide – Text: [20]	ence: DMA] DX: [[382, Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen e finish in compariso	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church 1930 above/belo	Location Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1902 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON Question: was the Reasoning Type: Answer: above Bounding Box: [] Supporting Evide – Text: [20] – Bounding Bo	ence: DMA] DX: [[382, Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen e finish in compariso	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Levanger Church Markabygda Church Okkenhaug Chapel Ytterøy Church Asen Church 1930 above/belo Dn 358, 577, 388]]	Location Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytteraya Åsen	Year built 1180 1893 1992 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON Question: was the Reasoning Type: Answer: above Bounding Box: [] Supporting Evide – Text: [20] – Bounding Bo	ence: DMA] ox: [[382, Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen e finish in compariso	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Bamberg Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church 1930 above/belo Dn 358, 577, 388]]	Vanger Location Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1902 1998 1887 1893 1890 1904	
• A.4.14	Bounding Box: [] Supporting Evide – Text: [DIPLO – Bounding Bo COMPARISON Question: was the Reasoning Type: Answer: above Bounding Box: [] Supporting Evide – Text: [20] – Bounding Bo	Parish Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøy Åsen e finish in compariso	222, 610,268]] Churches in Le Church name Alstadhaug Church Ekne Church Bamberg Church Markabygda Church Okkenhaug Chapel Ytterøy Church Åsen Church 1930 above/belo Dn 358, 577, 388]]	Location Alstadhaug Ekne Levanger Markabygd Okkenhaug Ytterøya Åsen	Year built 1180 1893 1902 1998 1887 1893 1890 1904	

