

Figure 2: Overview of different tasks in our benchmark: Eight key components illustrating the task inputs and outputs for table recognition, chart understanding, text recognition, diagram analysis, VQA, line detection, layout analysis, and PDF-to-Markdown conversion, complete with input/output examples for each task.

Domain/ Characteristics	EXAMS-V*	Camel- Bench	MIDAD†	KHATT	KITAB- Bench (Ours)
PDF to Markdown	×	×	×	×	✓
Layout Detection	×	×	×	×	✓
Line Detection	×	×	×	×	✓
Line Recognition	×	✓	×	×	✓
Table Recognition	×	×	×	×	✓
Image to Text	✓	✓	✓	✓	✓
Charts to JSON	×	×	×	×	✓
Diagram to Code	×	×	×	×	✓
VQA	×	✓	×	×	✓
Handwritten Samples	×	×	×	✓	✓
Open Source	✓	✓	×	×	✓
Total Samples (#)	823	3,004	29,435	5,000	8,809

Table 1: Comparison of Arabic OCR Benchmarks Across Different Domains. Benchmarks compared: LaraBench (Abdelali et al., 2023), CamelBench (Ghaboura et al., 2024), MIDAD (Bhatia et al., 2024), KHATT (Mahmoud et al., 2014), and KITAB-Bench (Ours). (\*: Only the Arabic samples are considered.) (†: The test set of the dataset is considered.)

processing challenges such as table parsing, font detection, and numeral recognition. Arabic benchmarks like CAMEL-Bench (Ghaboura et al., 2024) and LaraBench (Abdelali et al., 2023) evaluate large multimodal and language models, but they give limited attention to document understanding tasks. Consequently, there remains a need for a more comprehensive framework to systematically evaluate and compare Arabic OCR solutions. Our benchmark addresses these gaps by offering diverse document types and evaluation tasks to facilitate in-depth assessments of modern OCR systems.

We present KITAB-Bench, a comprehensive Arabic OCR benchmark spanning 9 domains and 36 sub-domains. Our framework evaluates layout detection (text blocks, tables, figures), multi-format

recognition (printed/handwritten text, charts, diagrams), and structured output generation (HTML tables, DataFrame charts, markdown). This enables rigorous assessment of both basic OCR capabilities and advanced document understanding tasks.

The contributions of this work include (1) A comprehensive Arabic OCR benchmark covering multiple document types and recognition tasks. (2) Detailed evaluation metrics for assessing performance across different document understanding challenges. We also propose CharTeX and CODM metric to evaluate chart extraction and diagram extraction respectively. (3) Baseline results for popular OCR systems and Vision Language Models (VLMs), highlighting current limitations and areas for improvement. (4) A standardized framework for comparing Arabic OCR systems, facilitating future research and development.

## 2 Related Work

The development of robust Optical Character Recognition (OCR) systems has been extensively studied across document layout analysis (Zhao et al., 2024; Shen et al., 2021; Paruchuri, 2024b; JaiedAI, 2020; Auer et al., 2024; Li et al., 2020), table detection (Li et al., 2019; Paliwal et al., 2019; Nassar et al., 2022; Li et al., 2021; Schreiber et al., 2017), and document understanding (Staar et al., 2018; Weber et al., 2023; Livathinos et al., 2021). While English OCR benefits from rich datasets like PubLayNet (Zhong et al., 2019b), DocBank (Li et al., 2020), M6Doc (Cheng et al.,

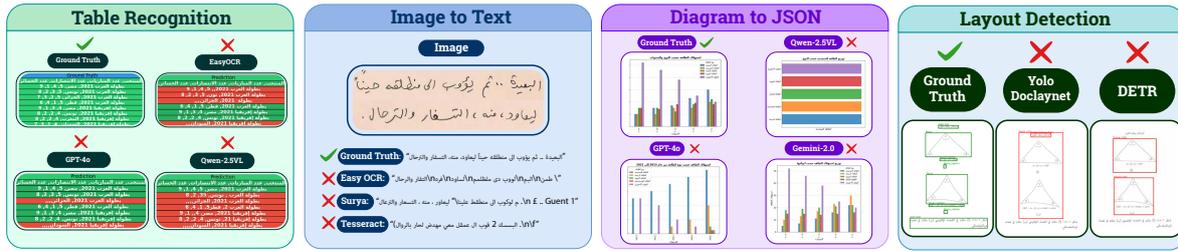


Figure 3: Comparison of model performance across four document understanding tasks (Table Recognition, Image to Text, Diagram to JSON, and Layout Detection) showing successful and failed cases for different models including Ground Truth, EasyOCR, GPT-4, Qwen, Surya, Tesseract, Yolo, and DETR on Arabic document benchmark data.

Domain	Total Samples
PDF to Markdown	33
Layout	2,100
Line Detection	378
Line Recognition	378
Table Recognition	456
Image to Text	3,760
Charts to DataFrame	576
Diagram to Json	226
VQA	902
<b>Total</b>	<b>8,809</b>

Table 2: Distribution of samples across different domains in our dataset. A more detailed count for different sub-domains and data sources is in Appendix A.

2023), and DocLayNet (Pfitzmann et al., 2022), Arabic lacks standardized benchmarks for diverse fonts and layouts. Recent efforts like MIDAD (Bhatia et al., 2024) curates extensive training data for Arabic OCR and handwriting recognition, while Peacock (Alwajih et al., 2024) introduces culturally-aware Arabic multimodal models. Existing resources such as CAMEL-Bench (Ghaboura et al., 2024), LARA-Bench (Abdelali et al., 2023), MADAR (Bouamor et al., 2018), OSACT (Mubarak et al., 2022), and Tashkeela (Zerrouki and Balla, 2017) focus on language modeling or specific tasks rather than full-page OCR evaluation. Handwriting datasets including HistoryAr (Pantke et al., 2014), IFN/ENIT (Pechwitz et al., 2002), KHATT (Mahmoud et al., 2014), APTI (Slimane et al., 2009), and Muharaf (Saeed et al., 2024) emphasize word/line recognition over document structure analysis. Arabic table recognition faces challenges from merged cells and RTL formatting (Pantke et al., 2014). While methods like GTE (Zheng et al., 2021), GFTE (Li et al., 2021), CascadeTabNet (Prasad et al., 2020), TableNet (Paliwal et al., 2019), and TableFormer (Nassar et al., 2022) advance Latin table detection, their effectiveness on Arabic

documents remains unproven. Document conversion pipelines (CCS (Staar et al., 2018), Tesseract (Smith, 2007), Docling (Auer et al., 2024), Surya (Paruchuri, 2024b), Marker (Paruchuri, 2024a), MinerU (Wang et al., 2024a), PaddleOCR (Du et al., 2020)) lack Arabic-specific optimizations for segmentation and diacritic handling (Mahmoud et al., 2018; Kiessling et al., 2019). This highlights the critical need for comprehensive Arabic OCR benchmarks addressing text recognition, table detection, and layout parsing.

### 3 KITAB-Bench

Our methodology offers a novel approach to benchmarking Arabic OCR systems via a comprehensive data collection strategy and a systematic evaluation framework. We gather curated samples from existing Arabic document datasets, manually collected and annotated PDFs, and employ a five-phase LLM-assisted human-in-the-loop pipeline (Figure 4) to generate diverse supplementary content. Our evaluation framework spans nine specialized tasks, enabling thorough assessment of OCR performance across various document processing challenges and providing a robust benchmark for Arabic document understanding tasks.

#### 3.1 PDF Data Collection

We curated 33 diverse PDFs from online sources in academia, medicine, law, and literature. To ensure challenging cases, we selected documents featuring richly formatted tables with extensive color usage, merged cells, Arabic numerals, historical texts, watermarks, and handwritten annotations. Each PDF averaged three pages, and we then manually annotated them. This dataset comprehensively captures real-world complexities, making it a valuable benchmark for PDF-to-Markdown conversion.

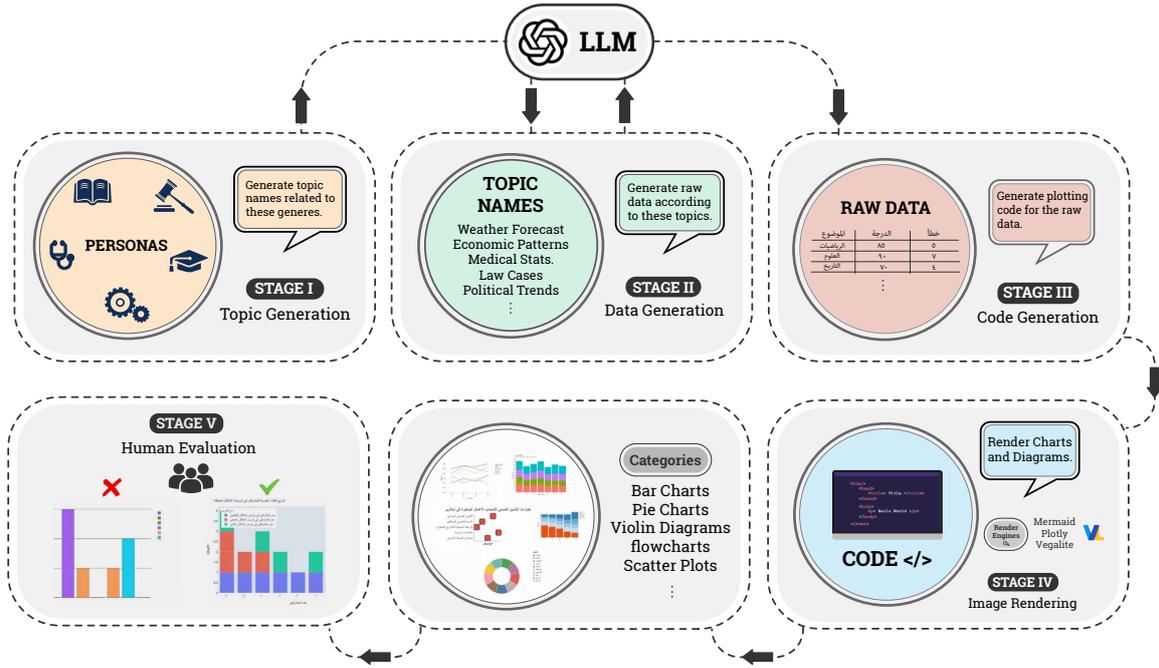


Figure 4: Synthetic Data Generation Pipeline: A 5-stage process using LLMs to generate topics, create raw data, produce visualization code, render charts, and perform human evaluation for quality control.

### 3.2 LLM-Assisted Data Generation Pipeline

To generate data for charts, diagrams and tables, we implemented a five-phase LLM-assisted generation pipeline with human validation at critical stages, as illustrated in Figure 4. *In Phase I (Topic Generation)*, our system employs an LLM to generate diverse topic names across multiple domains. This phase incorporates various personas (academic, legal, medical, technical) to ensure broad coverage of document types. *Phase II (Data Generation)* transforms the validated topics into structured raw data. The LLM generates content following Arabic linguistic and formatting conventions across various domains. *In Phase III (Code Generation)*, the system converts the validated raw data into plotting code, with special attention to Arabic text rendering requirements and RTL content management. *Phase IV (Image Rendering)* utilizes specialized rendering engines (Mermaid, Plotly, Vegalite, HTML) to create visual representations while maintaining Arabic text integrity.

*The final phase (Human Evaluation)* implements rigorous quality control through expert validation. Evaluators filter charts, tables and diagrams based on detected anomalies and ensure adherence to Arabic-specific document conventions. This phase is crucial for maintaining the high quality of our benchmark dataset.

### 3.3 Dataset Statistics

Our benchmark dataset comprises over 8,809 samples across 9 major domains and 36 sub-domains, representing a comprehensive collection of Arabic document types for OCR evaluation. As detailed in Table 8, the dataset combines carefully curated samples from established datasets, manually annotation PDFs, and synthetically generated content created through our LLM-assisted pipeline (Figure 4). The Image-to-Text portion (3,760 samples) includes data from historical documents (HistoryAr (Pantke et al., 2014)), handwritten text collections (Khatt (Mahmoud et al., 2014), ADAB (Boubaker et al., 2021), Muharaf (Saeed et al., 2024)), and scene text (EvAREST (Hassan et al., 2021)), while layout detection comprises 2,100 samples from BCE-Arabic-v1 (Saad et al., 2016) and DocLayNet (Pfitzmann et al., 2022).

For layout analysis, we incorporated 1,700 samples from BCE-Arabic-v1 dataset (Saad et al., 2016), 400 samples from DocLayNet dataset (Pfitzmann et al., 2022) focusing on financial, academic, legal, and patent documents. The line detection and recognition tasks contains 378 samples each from self-developed dataset. We further enriched the dataset with 500 samples from PATS-A01 (El-Muhtaseb, 2010) benchmark to ensure diverse representation. For handwritten text recognition, we assembled a comprehensive collection of 1,000

Task	Metric	Surya	Tesseract	EasyOCR
Detection	mAP@50	<b>79.67</b>	46.39	68.02
	mAP@0.5:0.95	27.40	14.30	<b>32.74</b>
Recognition	WER	1.01	1.00	<b>0.53</b>
	CER	0.87	0.66	<b>0.20</b>

Table 3: Performance of different models on Line Detection and Line Recognition Task on our Benchmark

samples combining datasets from Khatt (Mahmoud et al., 2014) (both paragraph and line-level annotations), Adab (Boubaker et al., 2021), Muharaf (Saeed et al., 2024), and OnlineKhatt (Mahmoud et al., 2018). The benchmark also includes specialized content from ISI-PPT (Wu and Natarajan, 2017) (500 samples), and Hindawi (Elfilali, 2023) (200 samples) for various document types. Scene text understanding is supported by 800 samples from EvArest (Hassan et al., 2021), providing real-world context diversity. A detailed table showing all the dataset is provided in the Appendix A.

A significant portion of our dataset consists of synthetically generated content, including 576 samples for Charts-to-DataFrame (spanning 16 different chart types), 422 samples for Diagram-to-Code (covering sequence diagrams, flowcharts, and tree maps), 456 samples for Tables-to-CSV/HTML, and 902 samples for VQA tasks. These synthetic samples were generated through our five-phase LLM-assisted human-in-the-loop pipeline (Figure 4). Every sample in our dataset - whether from existing sources or newly generated - underwent validation by native Arabic speakers before inclusion in the final benchmark. This rigorous validation, reinforced by expert review and automated checks, ensures high quality and authenticity across all domains. A detailed analysis is in Appendix C.

## 4 Experiments

Our experimental evaluation comprehensively assesses the capabilities of current OCR systems and state-of-the-art vision-language models (VLMs) across different Arabic and multilingual document understanding tasks. Figure 2 illustrates the nine distinct tasks in our evaluation framework.

We evaluate three categories of systems: VLMs, traditional OCR systems, and specialized document processing tools. For VLMs, we include both closed-source models like gpt-4o-2024-08-06, gpt-4o-mini-2024-07-18 (Hurst et al., 2024; Achiam et al., 2023), and gemini-2.0-flash (Georgiev et al., 2024; Google DeepMind,

2025), as well as open-source alternatives such as Qwen2-VL-7B (Wang et al., 2024b), Qwen2.5-VL-7B (Team, 2025), and the AIN-7B (Heakl et al., 2025). Traditional OCR approaches in our evaluation include Surya (Paruchuri, 2024b), Tesseract (Smith, 2007), EasyOCR (JaidedAI, 2020), and PaddleOCR (Li et al., 2022; Du et al., 2021). For specialized document processing tasks, we employ systems like Docling (Auer et al., 2024), and Marker (Paruchuri, 2024a). Layout detection capabilities are evaluated using methods implemented in Surya-layout (Paruchuri, 2024b), Yolo-doctraynet (Zhao et al., 2024) from MinerU (Wang et al., 2024a), and RT-DETR (Zhao et al., 2023) based method in Docling (Auer et al., 2024).

### 4.1 Evaluation Frameworks and Metrics

Our evaluation framework comprises nine specialized tasks designed to assess different aspects of Arabic OCR systems, as demonstrated in Figure 2. Each task addresses specific challenges in Arabic document processing. For this reason, we employ task-specific metrics to evaluate different aspects of document understanding.

**PDF-to-Markdown:** It evaluates the conversion of Arabic PDFs to structured markdown while preserving the text and table structure. Since both table and text structure are important, for evaluating PDF to Markdown conversion quality, we propose MARS (Markdown Recognition Score), which combines chrF (Popović, 2015) with Tree-Edit-Distance-based Similarity (TEDS) (Zhong et al., 2020) :

$$\text{MARS} = \alpha \cdot \text{chrF}_3 + (1 - \alpha) \cdot \text{TEDS}(T_a, T_b) \quad (1)$$

where  $\alpha$  ( $0 \leq \alpha \leq 1$ ) is the weight.  $T_a$  represent

Dataset	Metric	Surya	Yolo-doc-laynet	Detr (docling)
BCE	mAP@0.5	0.506	0.470	<b>0.750</b>
	mAP@0.5:0.95	0.381	0.369	<b>0.566</b>
	Precision	<b>0.751</b>	0.608	0.626
	Recall	0.593	0.592	<b>0.725</b>
	F1 Score	0.635	0.585	<b>0.654</b>
DocLayNet	mAP@0.5	0.675	0.404	<b>0.758</b>
	mAP@0.5:0.95	0.469	0.335	<b>0.541</b>
	Precision	<b>0.782</b>	0.527	0.635
	Recall	0.856	0.503	<b>0.770</b>
	F1 Score	0.799	0.499	<b>0.670</b>

Table 4: Performance comparison of layout detection models using different evaluation metrics

296 predicted table structure and  $T_b$  the ground truth  
297 structure.

**Table Recognition:** We evaluate table extraction using both HTML and CSV formats, where HTML format (evaluated using TEDS (Zhong et al., 2020)) preserves rich structural information including cell spans and hierarchical relationships crucial for complex Arabic tables, while CSV format (evaluated using Jaccard Index 2) focuses on raw data extraction optimized for machine processing and data analysis pipelines. This dual-format evaluation ensures systems can both maintain complex table structures for human readability and provide clean, structured data for automated processing, specifically important for RAG based systems.

$$J(P, G) = \frac{|P \cap G|}{|P \cup G|} = \frac{|P \cap G|}{|P| + |G| - |P \cap G|} \quad (2)$$

298 where  $|P \cap G|$  represents the number of exact  
299 matching cells between predicted and ground truth  
300 tables, and  $|P \cup G|$  represents the total number of  
301 unique cells across both tables.

**Chart-to-Dataframe:** This task evaluates extracting structured data from Arabic charts into machine-readable dataframes. Systems must accurately parse numerical values, text labels, and preserve data relationships across chart types (bar, line, pie). We use the Structuring Chart-oriented Representation Metric (SCRM) (Xia et al., 2024)—which combines type recognition, topic understanding, and structural numerical fidelity (see Appendix D)—and also propose our own CharTeX (Chart Extraction Score) metric. CharTeX combines the cHrf scores for chart type and topic with the jaccard index for the dataframe, using fuzzy matching (80% threshold) when columns do not exactly align.

$$\text{Metric} = \alpha J_{type} + \beta J_{topic} + (1 - \alpha - \beta) J_{data} \quad (3)$$

302 Here,  $J_{type}$  and  $J_{topic}$  denote the chrF scores be-  
303 tween the predicted and ground-truth chart type and  
304 topic, while  $J_{data}$  measures the structural similarity  
305 of the predicted and ground-truth JSON data.

306 **Diagram-to-JSON:** This task evaluates the con-  
307 version of Arabic flowcharts and technical diag-  
308 grams into JSON while preserving semantic rela-  
309 tionships and technical specifications. We propose  
310 CODM (Code-Oriented Diagram Metric), extend-  
311 ing SCRM (Xia et al., 2024), with the same fom-  
312 ulation as in Eq 3. More detail about this metric is  
313 provided in Appendix E

314 **Image-to-Text:** This task assess the basic text

315 recognition capabilities across different Arabic  
316 fonts and styles, including the handling of cursive  
317 script connections, diacritical marks, and various  
318 text orientations. We use we use Character Error  
319 Rate (CER) and Word Error Rate (WER). For a pre-  
320 dicted text sequence  $\hat{y}$  and ground truth sequence  $y$ ,  
321 CER is computed as:  $\text{CER} = \frac{L(y, \hat{y})}{|y|}$ , where  $L(y, \hat{y})$   
322 is the Levenshtein distance between character se-  
323 quences and  $|y|$  is the ground truth length. WER is  
324 calculated the same way with words as the unit of  
325 error.

**Visual Question Answering:** Tests the ability of  
326 models to understand and reason about Arabic doc-  
327 ument content, we evaluate using standard accuracy  
328 for MCQ questions and exact word match.  
329

**Line Detection:** Focuses on the accurate iden-  
330 tification and processing of individual text lines  
331 in Arabic documents. We evaluate using mean  
332 Average Precision (mAP) at different Intersec-  
333 tion over Union (IoU) thresholds: mAP@0.5 and  
334 mAP@0.5:0.95, which assess the localization ac-  
335 curacy of detected text lines.  
336

**Layout Detection:** Assesses document structure  
337 analysis capabilities, including the identification of  
338 headers, paragraphs, and complex layout elements  
339 in Arabic documents. Performance is measured us-  
340 ing mAP@0.5 and mAP@0.5:0.95 for localization  
341 accuracy, complemented by Precision, Recall, and  
342 F1 scores to evaluate the overall detection quality  
343 across different layout components.  
344

All metrics are computed on our diverse bench-  
345 mark dataset, which encompasses various docu-  
346 ment types and complexity levels in both Arabic  
347 and multilingual contexts. Table 10 provides a  
348 detailed mapping of tasks, metrics, and evaluated  
349 systems.  
350

## 351 4.2 Experimental Setup

352 We implement our evaluation pipeline with care-  
353 ful consideration of hyperparameters for different  
354 metric. All experiments use NVIDIA A100 GPUs.  
355 For VLMs, we use their official implementations  
356 or API endpoints. Traditional OCR systems are  
357 evaluated using pre-trained models provided by  
358 the frameworks. For PDF-to-Markdown evaluation  
359 metric MARS 1, we choose  $\alpha = 0.5$  and  $\alpha = 0.5$   
360 and  $\beta = 0.2$  for Diagram-to-JSON evaluation met-  
361 ric CODM. We average the results over multiple  
362 runs, with performance comparisons shown in dif-  
363 ferent tables [ 3, 6, 5, 7, and 4].

Model Group	Models	Table Extraction		End-to-End PDF		
		TEDS (HTML)	Jaccard (CSV)	CHrF (Text)	TEDS (Table)	MARS
Closed	GPT-4o	<b>85.76</b>	<b>66.36</b>	69.62	<b>60.61</b>	65.12
	GPT-4o-mini	69.32	49.50	56.59	52.69	54.64
	Gemini-2.0-Flash	83.08	65.55	<b>75.75</b>	55.55	<b>65.65</b>
Open	Qwen2-VL-7B	57.83	40.20	40.30	2.54	21.42
	Qwen2.5-VL-7B	59.31	59.58	69.21	11.65	40.43
	AIN-7B	75.94	64.83	56.52	49.32	52.92
Framework	Tesseract	28.23 <sup>D</sup>	14.85 <sup>D</sup>	59.91 <sup>D</sup>	45.44 <sup>D</sup>	52.68 <sup>D</sup>
		38.64 <sup>I</sup>	16.04 <sup>I</sup>			
	EasyOCR	49.10 <sup>D</sup>	23.83 <sup>D</sup>	57.46 <sup>D</sup>	51.12 <sup>D</sup>	54.29 <sup>D</sup>
		39.09 <sup>I</sup>	17.88 <sup>I</sup>			
	Surya	50.15 <sup>M</sup>	70.42 <sup>M</sup>	58.38 <sup>M</sup>	44.29 <sup>M</sup>	51.34 <sup>M</sup>

<sup>D</sup>Docling (Auer et al., 2024) pipeline    <sup>I</sup>Img2Table (Cattan, 2021) pipeline    <sup>M</sup>Marker (Paruchuri, 2024a) pipeline

Table 5: Performance comparison of different models for table extraction and end-to-end PDF to markdown conversion tasks on our benchmark.

Group	Models	CHrF ↑	CER ↓	WER ↓
Closed	GPT-4o	61.01	0.31	0.55
	GPT-4o-mini	47.21	0.43	0.71
	Gemini-2.0-Flash	77.95	<b>0.13</b>	0.32
Open	Qwen2VL-7B	33.94	1.48	1.55
	Qwen2.5VL-7B	49.23	1.20	1.41
	AIN-7B	<b>78.33</b>	0.20	<b>0.28</b>
Framework	Tesseract	39.62	0.54	0.84
	EasyOCR	45.47	0.58	0.89
	Paddle	16.73	0.79	1.02
	Surya	20.61	4.95	5.61

Table 6: Performance comparison of models for OCR (image to text) tasks on our benchmark. A detailed performance comparison among different open-source dataset is available in Appendix B

## 5 Results and Discussion

In this section, we present a comprehensive evaluation of different models across different tasks of our framework. The results provide a clear distinction between the performance of closed-source models, open-source models, and framework-based solutions, revealing both their strengths and limitations. We observe very clear performance gap between closed and open-source solutions. While closed-source models like Gemini-2.0-Flash consistently outperform other models almost all the tasks.

### 5.1 Charts, Diagrams, and VQA

Table [7] presents model performance across different chart and diagram understanding tasks, evaluated using SCRM and CharTeX (for charts), and VQA-based accuracy metrics. Among closed-source models, Gemini-2.0 achieves the highest performance on chart understanding metrics, scor-

ing 71.4% on SCRM and 56.28% on CharTeX. The performance gap between Gemini-2.0 and GPT-4o is particularly pronounced in CharTeX evaluation (10.33%) compared to SCRM (2.8%). Open-source models shows a significant limitation in complex chart understanding. While their SCRM scores remain competitive, both Qwen variants score below 23% on CharTeX evaluation. The visual question-answering results reveal an important exception to the general closed-source advantage. AIN achieves 87% on PATDVQA, surpassing Gemini-2.0 by 11.5%. AIN also shows competitive performance on MTVQA (31.50%), which is similar to GPT-4o and 4% better than GPT-4o-mini. This shows that open-source models can be competitive with closed-source alternatives.

### 5.2 Layout and Lines: Document Structure

Our evaluation of document structure understanding reveals distinct performance patterns across layout detection and line processing tasks. In layout detection (Table 4), RT-DETR (Zhao et al., 2023) achieves superior overall performance with mAP0.5 scores of 0.750 and 0.758 on BCE (arabic only) and DocLayNet (english) dataset respectively. However, Surya (Paruchuri, 2024b) demonstrates higher precision (0.782 on DocLayNet, 0.751 on BCE), despite lower recall rates. This trade-off suggests that different architectures optimize for different aspects of layout detection.

The line processing results (Table 3) highlight a clear contrast between detection and recognition capabilities. While Surya excels in detection with a mAP@0.50 of 79.67%, EasyOCR demonstrates superior recognition performance (WER: 0.53, CER:

Group	Model	Chart		Diagram	Visual QA				Average
		SCRM	CharTeX	CODM	MTVQA <sup>O</sup>	ChartsVQA <sup>M</sup>	DiagramsVQA <sup>M</sup>	PATDVQA <sup>M</sup>	
Closed	GPT-4o	68.6	45.95	61.6	32.00	77.00	85.29	82.50	69.19
	GPT-4o-mini	67.2	43.33	61.4	26.80	58.00	83.33	80.00	62.03
	Gemini-2.0-Flash	<b>71.4</b>	<b>56.28</b>	<b>71.8</b>	<b>35.00</b>	72.00	<b>88.24</b>	75.50	67.68
Open	Qwen2-VL-7B	56.6	21.59	63.0	19.60	59.00	82.35	77.50	59.61
	Qwen2.5-VL-7B	36.2	22.08	59.2	23.00	74.00	79.41	74.50	62.72
	AIN-7B	66.6	34.61	66.40	31.50	75.00	85.29	<b>87.00</b>	<b>69.69</b>

Table 7: Model Performance on Chart Understanding, Diagram Parsing, and Visual Question Answering Tasks. For VQA tasks, *O* denotes open-ended question type from MTVQA (Tang et al., 2024) dataset and *M* denotes MCQ type questions.

0.20). This inverse relationship between detection and recognition performance across models indicates a fundamental challenge in optimizing both capabilities simultaneously. Notably, Tesseract shows consistent but lower performance across both metrics, suggesting that newer architectures have made significant improvements over traditional approaches. We also observe that no single model excels at both detection and recognition, which requires for hybrid solutions.

### 5.3 Tables, OCR, and PDF-to-Markdown

Across table extraction tasks (Table 5), closed-source models maintain a clear advantage, with GPT-4o achieving 85.76% TEDS and 66.36% Jaccard scores. Among open-source models, AIN (75.94% TEDS) significantly outperforms Qwen variants, while specialized frameworks like Surya achieve competitive results (70.42% Jaccard) through targeted pipelines. In OCR evaluation (Table 6), Gemini-2.0-Flash leads with the lowest error rates (CER: 0.13, WER: 0.32). Notably, AIN matches this performance level (WER: 0.28), while traditional OCR frameworks like EasyOCR and Tesseract show moderate performance (CER: 0.58, 0.54). The significant performance drop in Paddle (CER: 0.79) and Surya (CER: 4.95) highlights the challenges in developing robust OCR systems.

End-to-end document processing (Table 5) reveals the largest gaps between approaches. Closed-source models maintain consistent performance (GPT-4o: 65.12% MARS, Gemini-2.0: 65.65% MARS), while open-source models show substantial degradation (Qwen2-VL-7B: 21.42% MARS). Framework approaches achieve better stability, with Tesseract and EasyOCR scoring above 50% MARS, suggesting that specialized pipelines can partially bridge the gap with larger models in complete document processing tasks.

Our comprehensive evaluation demonstrates that while closed-source models maintain superior performance over open-source models across most tasks, specialized frameworks like Surya, RT-DETR Layout, and EasyOCR achieve competitive performance in targeted scenarios like table extraction, layout detection, and text recognition respectively. However, this framework advantage significantly diminishes in end-to-end pdf-to-markdown tasks where the integration capabilities of large models prove crucial, as evidenced by the performance gaps between commercial VLMs and traditional systems like EasyOCR, Surya and Tesseract in End-to-End PDF task (Table 5).

## 6 Conclusion

We introduce a comprehensive benchmark for Arabic OCR that fills the gap in standardized evaluation frameworks for Arabic document processing. Our dataset of 8,809 samples across nine major domains is the most diverse collection assembled for OCR evaluation, incorporating handwritten, scanned, synthetic, and scene text, as well as complex tables, charts, and end-to-end pdf-to-markdown. This framework extends beyond simple text recognition to include structural document analysis and enables systematic assessment of OCR performance across various fonts, styles, and layouts.

## 7 Limitations and Future Directions

Despite its contributions, this benchmark has limitations. While it covers diverse Arabic document types, it lacks full representation of historical manuscripts and low-resource dialects. Future work should expand to include these, along with scanned records from government, academic, and financial institutions.

490	Another key limitation is in table and chart recog-	Xavier Cattani. 2021. img2table: Extract tables from	542
491	nition, where OCR models struggle with struc-	images and scanned pdfs. <a href="https://github.com/xavctn/img2table">https://github.com/xavctn/img2table</a> . Accessed: 2025-02-14.	543
492	ture preservation, header detection, and merged		544
493	cell parsing. Though our benchmark introduces	H. Cheng, P. Zhang, S. Wu, et al. 2023. M6doc:	545
494	challenges in these areas, further refinements are	A large-scale multi-format, multi-type, multi-layout,	546
495	needed for robust multimodal OCR capable of	multi-language, multi-annotation category dataset for	547
496	jointly processing text, tables, and figures. Fu-	modern document layout analysis. In <i>Proceedings of the</i>	548
497	ture advancements should focus on dataset expan-	<i>IEEE/CVF Conference on Computer Vision and Pattern</i>	549
498	sion, novel evaluation metrics, deep learning re-	<i>Recognition (CVPR)</i> .	550
499	finements, and cross-lingual OCR innovations to	Matt Deitke, Christopher Clark, Sangho Lee, Rohun	551
500	enhance Arabic VLMs.	Tripathi, Yue Yang, Jae Sung Park, Mohammadreza	552
		Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini,	553
		et al. 2024. Molmo and pixmo: Open weights and	554
		open data for state-of-the-art multimodal models. <i>arXiv</i>	555
501	<b>References</b>	<i>preprint arXiv:2409.17146</i> .	556
502	Ahmed Abdelali, Hamdy Mubarak, Shammur Absar	Yuning Du, Chenxia Li, Ruoyu Guo, Cheng Cui, Wei-	557
503	Chowdhury, Maram Hasanain, Basel Mousi, Sabri	wei Liu, Jun Zhou, Bin Lu, Yehua Yang, Qiwen	558
504	Boughorbel, Yassine El Kheir, Daniel Izham, Fahim	Liu, Xiaoguang Hu, et al. 2021. Pp-ocrv2: Bag of	559
505	Dalvi, Majd Hawasly, et al. 2023. Larabench: Bench-	tricks for ultra lightweight ocr system. <i>arXiv preprint</i>	560
506	marking arabic ai with large language models. <i>arXiv</i>	<i>arXiv:2109.03144</i> .	561
507	<i>preprint arXiv:2305.14982</i> .		
508	Josh Achiam, Steven Adler, Sandhini Agarwal, Lama	Yuning Du, Chenxia Li, Ruoyu Guo, Xiaoting Yin,	562
509	Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo	Weiwei Liu, Jun Zhou, Yifan Bai, Zilin Yu, Yehua	563
510	Almeida, Janko Altschmidt, Sam Altman, Shyamal	Yang, Qingqing Dang, et al. 2020. Pp-ocr: A prac-	564
511	Anadkat, et al. 2023. Gpt-4 technical report. <i>arXiv</i>	tical ultra lightweight ocr system. <i>arXiv preprint</i>	565
512	<i>preprint arXiv:2303.08774</i> .	<i>arXiv:2009.09941</i> .	566
513	Fakhraddin Alwajih, El Moatez Billah Nagoudi, Gagan	Husni A. El-Muhtaseb. 2010. Pats-a01 - an ara-	567
514	Bhatia, Abdelrahman Mohamed, and Muhammad	bic text database. <a href="https://faculty.kfupm.edu.sa/ics/muhtaseb/ArabicOCR/PATS-A01.htm">https://faculty.kfupm.edu.sa/ics/muhtaseb/ArabicOCR/PATS-A01.htm</a> . Database	568
515	Abdul-Mageed. 2024. Peacock: A family of arabic mul-	for Arabic Text Recognition Research.	569
516	timodal large language models and benchmarks. <i>arXiv</i>		570
517	<i>preprint arXiv:2403.01031</i> .	Ali Elfilali. 2023. Hindawi books dataset.	571
518	Christoph Auer, Maksym Lysak, Ahmed Nassar,	<a href="https://huggingface.co/datasets/Ali-C137/Hindawi-Books-dataset">https://huggingface.co/datasets/Ali-C137/Hindawi-Books-dataset</a> . Dataset.	572
519	Michele Dolfi, Nikolaos Livathinos, Panos Vagenas,		573
520	Cesar Berrospi Ramis, Matteo Omenetti, Fabian Lindl-	Ling Fu, Biao Yang, Zhebin Kuang, Jiajun Song, Yuzhe	574
521	bauer, Kasper Dinkla, et al. 2024. Docling technical	Li, Linghao Zhu, Qidi Luo, Xinyu Wang, Hao Lu,	575
522	report. <i>arXiv preprint arXiv:2408.09869</i> .	Mingxin Huang, et al. 2024. Ocrbench v2: An improved	576
523	Gagan Bhatia, El Moatez Billah Nagoudi, Fakhraddin	benchmark for evaluating large multimodal models on	577
524	Alwajih, and Muhammad Abdul-Mageed. 2024. Qalam:	visual text localization and reasoning. <i>arXiv preprint</i>	578
525	A multimodal llm for arabic optical character and hand-	<i>arXiv:2501.00321</i> .	579
526	writing recognition. <i>arXiv preprint arXiv:2407.13559</i> .	Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai,	580
527	Houda Bouamor, Nizar Habash, Mohammad Salameh,	Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng	581
528	Wajdi Zaghrouani, Owen Rambow, Dana Abdulrahim,	Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking	582
529	Ossama Obeid, Salam Khalifa, Fadhl Eryani, Alexander	multimodal understanding across millions of tokens of	583
530	Erdmann, et al. 2018. The madar arabic dialect corpus	context. <i>arXiv preprint arXiv:2403.05530</i> .	584
531	and lexicon. In <i>Proceedings of the eleventh interna-</i>	Sara Ghaboura, Ahmed Heakl, Omkar Thawakar, Ali	585
532	<i>tional conference on language resources and evaluation</i>	Alharthi, Ines Riahi, Abduljalil Saif, Jorma Laaksonen,	586
533	<i>(LREC 2018)</i> .	Fahad S Khan, Salman Khan, and Rao M Anwer. 2024.	587
534	Houcine Boubaker, Abdelkarim Elbaati, Najiba	Camel-bench: A comprehensive arabic lmm benchmark.	588
535	Tagougui, Haikal El Abed, Monji Kherallah, Volker	<i>arXiv preprint arXiv:2410.18976</i> .	589
536	Märgner, and Adel M. Alimi. 2021. <i>Adab database</i> .	Google DeepMind. 2025. <i>Gemini Model Updates -</i>	590
537	Hassina Bouressace and Janos Csirik. 2019. Printed	<i>February 2025</i> . Accessed: 2025-02-14.	591
538	arabic text database for automatic recognition systems.	Heba Hassan, Ahmed El-Mahdy, and Mohamed E Hus-	592
539	In <i>Proceedings of the 2019 5th International Conference</i>	sein. 2021. Arabic scene text recognition in the deep	593
540	<i>on Computer and Technology Applications</i> , pages 107–	learning era: Analysis on a novel dataset. <i>IEEE Access</i> ,	594
541	111.	9:107046–107058.	595



707	M. Saeed, A. Chan, A. Mijar, and J. Moukarzel. 2024. Muharaf: Manuscripts of handwritten arabic dataset for cursive text recognition. <i>arXiv preprint arXiv:2406.09630</i> .	763
708		764
709		765
710		766
711	Sebastian Schreiber, Stefan Agne, Ivo Wolf, Andreas Dengel, and Sheraz Ahmed. 2017. Deepdesrt: Deep learning for detection and structure recognition of tables in document images. In <i>2017 14th IAPR international conference on document analysis and recognition (ICDAR)</i> , volume 1, pages 1162–1167. IEEE.	767
712		768
713		769
714		770
715		771
716		772
717	Zejiang Shen, Ruochen Zhang, Melissa Dell, Benjamin Charles Germain Lee, Jacob Carlson, and Weining Li. 2021. Layoutparser: A unified toolkit for deep learning based document image analysis. In <i>Document Analysis and Recognition-ICDAR 2021: 16th International Conference, Lausanne, Switzerland, September 5–10, 2021, Proceedings, Part I 16</i> , pages 131–146. Springer.	773
718		774
719		775
720		776
721		777
722		778
723		779
724	Fouad Slimane, Rolf Ingold, Slim Kanoun, Adel M Alimi, and Jean Hennebert. 2009. A new arabic printed text image database and evaluation protocols. In <i>2009 10th international conference on document analysis and recognition</i> , pages 946–950. IEEE.	780
725		781
726		782
727		783
728		784
729	Ray Smith. 2007. An overview of the tesseract ocr engine. In <i>Ninth international conference on document analysis and recognition (ICDAR 2007)</i> , volume 2, pages 629–633. IEEE.	785
730		786
731		787
732		788
733	Peter WJ Staar, Michele Dolfi, Christoph Auer, and Costas Bekas. 2018. Corpus conversion service: A machine learning platform to ingest documents at scale. In <i>Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery &amp; Data Mining</i> , pages 774–782.	789
734		790
735		791
736		792
737		793
738		794
739	Jingqun Tang, Qi Liu, Yongjie Ye, Jinghui Lu, Shu Wei, Chunhui Lin, Wanqing Li, Mohamad Fitri Faiz Bin Mahmood, Hao Feng, Zhen Zhao, Yanjie Wang, Yuliang Liu, Hao Liu, Xiang Bai, and Can Huang. 2024. Mtvqa: Benchmarking multilingual text-centric visual question answering. <i>Preprint</i> , arXiv:2405.11985.	795
740		796
741		797
742		798
743		799
744		800
745	Qwen Team. 2025. <a href="#">Qwen2.5-vl</a> .	801
746		802
747	Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, et al. 2024a. Mineru: An open-source solution for precise document content extraction. <i>arXiv preprint arXiv:2409.18839</i> .	803
748		804
749		805
750		806
751	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024b. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .	807
752		808
753		809
754		810
755		811
756	Maurice Weber, Carlo Siebenshuh, Rory Butler, Anton Alexandrov, Valdemar Thanner, Georgios Tsolakis, Haris Jabbar, Ian Foster, Bo Li, Rick Stevens, et al. 2023. Wordscape: a pipeline to extract multilingual, visually rich documents with layout annotations from web crawl data. <i>Advances in Neural Information Processing Systems</i> , 36:26048–26068.	812
757		813
758		814
759		815
760		816
761		
762		
	Haoran Wei, Chenglong Liu, Jinyue Chen, Jia Wang, Lingyu Kong, Yanming Xu, Zheng Ge, Liang Zhao, Jianjian Sun, Yuang Peng, et al. 2024. General ocr theory: Towards ocr-2.0 via a unified end-to-end model. <i>arXiv preprint arXiv:2409.01704</i> .	
	Yue Wu and Prem Natarajan. 2017. Self-organized text detection with minimal post-processing via border learning. In <i>International Conference on Computer Vision</i> .	
	Renqiu Xia, Bo Zhang, Haoyang Peng, Hancheng Ye, Xiangchao Yan, Peng Ye, Botian Shi, Yu Qiao, and Junchi Yan. 2023. Structchart: Perception, structuring, reasoning for visual chart understanding. <i>arXiv preprint arXiv:2309.11268</i> .	
	Renqiu Xia, Bo Zhang, Hancheng Ye, Xiangchao Yan, Qi Liu, Hongbin Zhou, Zijun Chen, Min Dou, Botian Shi, Junchi Yan, et al. 2024. Chartx & chartvlm: A versatile benchmark and foundation model for complicated chart reasoning. <i>arXiv preprint arXiv:2402.12185</i> .	
	Taha Zerrouki and Amar Balla. 2017. Tashkeela: Novel corpus of arabic vocalized texts, data for auto-diacritization systems. <i>Data in brief</i> , 11:147.	
	Y Zhao, W Lv, S Xu, J Wei, G Wang, Q Dang, Y Liu, and J Chen. 2023. Detrs beat yolos on real-time object detection. arxiv e-prints. <i>arXiv preprint arXiv:2304.08069</i> .	
	Z. Zhao, H. Kang, B. Wang, and C. He. 2024. Doclayout-yolo: Enhancing document layout analysis through diverse synthetic data and global-to-local adaptive perception. <i>arXiv preprint arXiv:2410.12628</i> .	
	Xinyi Zheng, Douglas Burdick, Lucian Popa, Xu Zhong, and Nancy Xin Ru Wang. 2021. Global table extractor (gte): A framework for joint table identification and cell structure recognition using visual context. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 697–706.	
	X Zhong, E ShafieiBavani, and A Jimeno-Yepes. 2019a. Image-based table recognition: data, model, and evaluation. corr abs/1911.10683. <i>arXiv preprint arXiv:1911.10683</i> .	
	Xu Zhong, Elaheh ShafieiBavani, and Antonio Jimeno Yepes. 2020. Image-based table recognition: data, model, and evaluation. In <i>European conference on computer vision</i> , pages 564–580. Springer.	
	Xu Zhong, Jianbin Tang, and Antonio Jimeno Yepes. 2019b. Publaynet: largest dataset ever for document layout analysis. In <i>2019 International conference on document analysis and recognition (ICDAR)</i> , pages 1015–1022. IEEE.	

## A Source of the Existing Dataset Collection

Our benchmark integrates diverse data sources to ensure comprehensive coverage of Arabic document types. As detailed in Table 2, the dataset

817 combines manually curated samples, synthetic data  
 818 generated through our LLM-assisted pipeline (Fig-  
 819 ure 4), and existing publicly available datasets. Key  
 820 sources include:

- 821 • Handwritten Text: KHATT (paragraph and  
 822 line-level annotations), ADAB, Muharaf, and  
 823 OnlineKhatt.
- 824 • Historical Documents: HistoryAr and Histori-  
 825 calBooks.
- 826 • Scene Text: EvAREST for real-world context  
 827 diversity.
- 828 • Layout Analysis: BCE-Arabic-v1 and Do-  
 829 cLayNet.
- 830 • Synthetic Content: 576 chart samples (16  
 831 types) and 422 diagram samples generated  
 832 via our five-phase pipeline (Section 3.2).

833 The dataset emphasizes domain diversity, covering  
 834 academic, medical, legal, financial, and technical  
 835 documents. All samples underwent rigorous valida-  
 836 tion by native Arabic speakers to ensure linguistic  
 837 and structural accuracy.

## 838 B Detailed Performance Comparison

839 Table 9 provides granular performance metrics for  
 840 VLMs and OCR frameworks across 12 Arabic text  
 841 recognition datasets. Gemini-2.0-Flash demon-  
 842 strates exceptional robustness on synthetic datasets  
 843 (CER: 0.01 on PATS), while AIN-7B excels in his-  
 844 torical manuscript recognition (CER: 0.26 on His-  
 845 toryAr). Traditional OCR systems like Tesseract  
 846 show limitations in handwritten text (CER: 1.26 on  
 847 HistoryAr), highlighting the need for script-specific  
 848 optimizations.

## 849 C Data Analysis

850 Our data generation pipeline (Figure 4) enabled  
 851 the creation of 1,502 synthetic samples (576 charts,  
 852 422 diagrams, 456 tables). The pipeline’s human  
 853 validation phase rejected 18% of initial outputs  
 854 due to RTL formatting errors or semantic incon-  
 855 sistencies. As shown in Figure 5 and 6, domain-  
 856 specific prompts ensured adherence to Arabic lin-  
 857 guistic conventions during LLM-assisted genera-  
 858 tion. The final dataset exhibits balanced represen-  
 859 tation across:

- 860 • Font Styles: 21 Arabic calligraphic styles

- Document Types: 36 sub-domains including  
 financial reports and technical manuals 861 862
- Structural Complexity: 43% of tables contain  
 merged cells; 29% of charts use dual-axis con-  
 figurations 863 864 865

## D Tasks Models and Metrics 866

867 Table 10 maps evaluation tasks to corresponding  
 868 models and metrics. The framework evaluates nine  
 869 core capabilities:

- Structural Understanding: Layout detection  
 (mAP), line detection (IoU) 870 871
- Content Extraction: Text recognition (CER),  
 table parsing (TEDS) 872 873
- Semantic Reasoning: VQA accuracy, chart-  
 to-dataframe conversion (SCRM) 874 875
- Specialized metrics like MARS ( $\alpha=0.5$ ) ad-  
 dress the dual requirements of text fidelity and  
 structural preservation in PDF-to-Markdown  
 conversion. 876 877 878 879

## E SCRM and CODM 880

881 The Structuring Chart-oriented Representation Met-  
 882 ric (SCRM) evaluates chart understanding through  
 three components:

$$883 \text{SCRM} = 0.4J_{\text{type}} + 0.3J_{\text{topic}} + 0.3J_{\text{data}} \quad (4)$$

884 where  $J_{\text{type}}$ , and  $J_{\text{topic}}$  are chrF scores, and  $J_{\text{data}}$   
 885 measures JSON structural similarity.

The Code-Oriented Diagram Metric (CODM) ex-  
 tends SCRM for flowcharts and technical diagrams:

$$886 \text{CODM} = 0.5J_{\text{topology}} + 0.5J_{\text{semantics}} \quad (5)$$

887 assessing both node-edge relationships and seman-  
 888 tic labels. As shown in Figure 5 and 6, domain-  
 889 specific prompts guided model responses for metric  
 calculation. For instance, sequence diagrams re-  
 quired strict adherence to Arabic UML notation  
 standards during evaluation. 884 885 886 887 888

Domain	Sub-Domain	Dataset Source	Original	Selected	Total													
PDF to Markdown	General	Manual	33	33	33													
Layout Detection	Docs	BCE-Arabic-v1 (Saad et al., 2016) DocLayNet (Pfitzmann et al., 2022)	1.9k 80k	1,700 400	2,100													
Line Detection	Docs	Manual	375	378	378													
Line Recognition	Docs	Manual	375	378	378													
Table Recognition	Financial	Pixmo (Deitke et al., 2024)	490	456	456													
Image to Text	Synthetic	PATS (El-Muhtaseb, 2010)	21.6k	500	3,760													
		SythenAR	39.1k	500														
	Historical	HistoryAr (Pantke et al., 2014)	1.5k	200														
		HistoricalBooks	40	10														
	Hand. Paragraph	Khatt (Mahmoud et al., 2014)	2.72k	200														
		Hand. Word	ADAB (Boubaker et al., 2021)	15k		200												
	Hand. Line		Muharaf (Saeed et al., 2024)	24.5k		200												
		PPT	OnlineKhatt (Mahmoud et al., 2018)	8.5k		200												
	Blogs		Khatt (Mahmoud et al., 2014)	13.4k		200												
		Scene	ISI-PPT (Wu and Natarajan, 2017)	86.5k		500												
	Charts to DataFrame		ArabicOCR	20.3k		50												
		Charts to DataFrame	Hindawi (Elfilali, 2023)	79k		200												
	Charts to DataFrame		EvAREST (Hassan et al., 2021)	5.59k		800												
Charts to DataFrame			Bar	Synthetic	100	61												
			Charts to DataFrame	Line	Synthetic	100	43											
				Charts to DataFrame	Pie	Synthetic	100	56										
					Charts to DataFrame	Box	Synthetic	100	31									
						Charts to DataFrame	Violin	Synthetic	100	36								
							Charts to DataFrame	Area	Synthetic	50	29							
								Charts to DataFrame	SunBurst	Synthetic	30	15						
									Charts to DataFrame	Dot	Synthetic	30	15					
										Charts to DataFrame	Dual Axis	Synthetic	20	26				
											Charts to DataFrame	Density Curve	Synthetic	10	5			
												Charts to DataFrame	Bubble	Synthetic	20	13		
													Charts to DataFrame	Grouped Bar	Synthetic	50	60	
														Charts to DataFrame	Stacked Bar	Synthetic	50	82
															Charts to DataFrame	Histogram	Synthetic	100
		Charts to DataFrame														HeatMap	Synthetic	10
	Charts to DataFrame															Scatter	Synthetic	100
Diagram to Json																Sequence	Synthetic	50
			Diagram to Json													Funnel	Synthetic	20
				Diagram to Json												Class	Synthetic	20
					Diagram to Json											Network	Synthetic	20
						Diagram to Json										Venn	Synthetic	20
							Diagram to Json									FlowChart	Synthetic	100
								Diagram to Json								TreeMap	Synthetic	100
VQA									Diagrams							Manual	102	102
			VQA						Charts	Manual						105	100	
				VQA					News Letter	PATD (Bouessace and Csirik, 2019)	2.42k					200		
					VQA				Scene	MTVQA	818	500						
<b>Total Dataset Size</b>						-			8,809									

Table 8: Dataset Distribution Across Different Domains, sub-domains and Data Source

Dataset	Size	GPT-4o		GPT-4o-mini		Gemini-2.0-Flash		Qwen2-VL	
		CER	WER	CER	WER	CER	WER	CER	WER
PATS	500	0.23	0.30	0.53	0.71	0.01	0.02	1.02	1.02
SythenAR	500	0.09	0.20	0.14	0.32	0.07	0.17	0.59	1.13
HistoryAr	200	0.51	0.82	0.67	0.96	0.28	0.64	3.46	2.86
HistoricalBooks	10	0.41	0.76	0.59	0.88	0.05	0.22	1.90	2.16
Khatt	200	0.45	0.74	0.64	0.91	0.19	0.45	1.12	5.04
Adab	200	0.30	0.73	0.35	0.83	0.19	0.56	0.63	1.08
Muharaf	200	0.56	0.90	0.63	0.94	0.33	0.69	3.57	2.87
OnlineKhatt	200	0.29	0.63	0.41	0.76	0.17	0.44	1.30	2.01
ISI-PPT	500	0.08	0.18	0.15	0.31	0.06	0.15	1.03	1.06
ArabicOCR	50	0.06	0.26	0.16	0.46	0.00	0.02	1.25	1.50
Hindawi	200	0.34	0.56	0.48	0.71	0.01	0.04	1.82	2.05
EvArest	800	0.20	0.38	0.25	0.51	0.18	0.36	0.41	0.95
	3,760	0.31	0.55	0.43	0.71	0.13	0.32	1.48	1.20

Dataset	Size	Qwen2.5-VL		AIN		Tesseract		Surya	
		CER	WER	CER	WER	CER	WER	CER	WER
PATS	500	0.26	0.36	0.00	0.00	0.14	0.28	4.66	4.67
SythenAR	500	0.21	0.40	0.04	0.16	0.31	0.72	4.82	7.90
HistoryAr	200	0.47	0.83	0.26	0.54	0.72	1.26	10.32	12.78
HistoricalBooks	10	0.33	0.72	0.84	0.88	0.74	0.99	6.81	6.30
Khatt	200	0.07	0.22	0.61	1.12	0.67	1.06	4.25	3.77
Adab	200	0.00	0.01	1.00	1.00	1.00	1.14	7.28	8.71
Muharaf	200	0.61	0.96	0.38	0.54	0.77	1.22	6.19	7.48
OnlineKhatt	200	0.36	0.70	0.03	0.12	0.59	1.20	6.71	6.95
ISI-PPT	500	0.36	0.54	0.52	0.53	0.31	0.64	4.25	3.77
ArabicOCR	50	1.00	1.00	0.01	0.01	0.01	0.01	2.75	3.58
Hindawi	200	1.00	1.00	0.11	0.15	0.31	0.72	0.15	0.20
EvArest	800	0.19	0.36	0.30	0.32	0.85	1.02	5.91	3.86
	3,760	0.28	0.54	0.20	0.58	0.89	0.79	4.95	5.61

Table 9: Performance comparison of Large Vision-Language Models on KITAB-Bench (lower is better).

Task	Metrics	Open LLMs	Closed LLMs	OCR Systems
<i>Document Understanding Tasks</i>				
PDF to Markdown	chrF + TEDS	–	–	Docling Marker MinerU PDF-Extract-Kit
Layout Detection	mAP@0.5 mAP@0.5:0.95 Precision Recall F1	–	–	Surya Yolo-doclaynet (MinerU) Detr (docling)
Line Detection	mAP@0.5 mAP@0.5:0.95	–	–	Surya Tesseract EasyOCR
Line Recognition	WER, CER	–	–	Surya Tesseract EasyOCR
<i>Table Understanding Tasks</i>				
Tables Recognition (HTML)	TEDS (Zhong et al., 2019a)	Qwen2-VL Qwen2.5-VL AIN PaliGemma	GPT-4o GPT-4o-mini Gemini-2.0-Flash	Docling[EasyOCR] Docling[Tesseract] Marker Img2Table[EasyOCR] Img2Table[Tesseract]
Tables Recognition (CSV)	Jaccard Index	Qwen2-VL Qwen2.5-VL AIN PaliGemma	GPT-4o GPT-4o-mini Gemini-2.0-Flash	Docling[EasyOCR] Docling[Tesseract] Marker Img2Table[EasyOCR] Img2Table[Tesseract]
<i>Visual Understanding Tasks</i>				
Image to Text	CER, WER chrF, BLEU METEOR	Qwen2-VL Qwen2.5-VL AIN-7B PaliGemma	GPT-4o GPT-4o-mini Gemini-2.0-Flash	Docling[EasyOCR] Docling[Tesseract] Marker Img2Table[EasyOCR] Img2Table[Tesseract]
Charts to DataFrame	SCRM (Xia et al., 2024, 2023)	Qwen2-VL Qwen2.5-VL AIN PaliGemma	GPT-4o GPT-4o-mini Gemini-2.0-Flash	–
Diagram to Json	SCRM	Qwen2-VL Qwen2.5-VL AIN-7B PaliGemma	GPT-4o GPT-4o-mini Gemini-2.0-Flash	–
VQA	Accuracy + Word Match Score	Qwen2-VL Qwen2.5-VL AIN-7b PaliGemma	GPT-4o GPT-4o-mini Gemini-2.0-Flash	–

Table 10: Comprehensive evaluation metrics and models for document understanding tasks. The table is organized into three main categories: document understanding, table understanding, and visual understanding tasks. Each task is evaluated using specific metrics and implemented across various models and OCR systems.

### Charts: Type Prompt

""You are an expert in detecting chart types. Below are examples of the expected output format:

Example 1:  
bar chart

Example 2:  
scatter chart

Example 3:  
histogram

Your task is to determine the type of chart shown in the given image.

**\*\*Instructions:\*\***

- **\*\*Respond with only the chart type\*\*** (e.g., 'bar chart', 'scatter chart').
- **\*\*Do not include any additional text, explanations, or descriptions.\*\***
- **\*\*Ensure the output matches the format in the examples exactly.\*\***

Provide only the chart type in **\*\*single quotes\*\*** as shown in the examples above.

What type of chart is shown in the image? Don't output any extra text""

### PDF to Markdown Prompt

""Extract the text from the document in Markdown format, and extract the tables in HTML format.  
Do not add style or anything, just the text. Do not ever generate tables in markdown format. Give me the output, nothing else.""

### OCR Prompt

""Extract the text in the image. Give me the final text, nothing else.""

### Diagrams: Type Prompt

""You are an expert in detecting chart types. Below are examples of the expected output format:

Example 1:  
treemap

Example 2:  
flowchart

Example 3:  
diagram

Your task is to determine the type of chart shown in the given image.

**\*\*Instructions:\*\***

- **\*\*Respond with only the chart type\*\*** (e.g., 'flowchart', 'sequence').
- **\*\*Do not provide any explanations, descriptions, or additional text.\*\***
- **\*\*Ensure the output strictly follows the format shown in the examples.\*\***

What type of chart is shown in the image?""

### Charts: Topic Prompt

""أنت خبير في تحليل وتقييم المخططات البيانية. فيما يلي أمثلة توضح تنسيق الإجابة المتوقع:"

**\*\*1 مثال:\*\***

توزيع الكتب الأكثر مبيعاً حسب النوع الأدبي

**\*\*2 مثال:\*\***

آراء العملاء حول الموضوعات المثيرة للجدل في الكتب

**\*\*التعليمات:\*\***

- **\*\*حدد موضوع أو محتوى المخطط البياني فقط فقط.\*\***
- **\*\*اكتب الإجابة باللغة العربية فقط.\*\***
- **\*\*اتبح التنسيق المحدد دون إضافة أي شرح أو تعليق إضافي.\*\***

""ما هو موضوع أو محتوى المخطط البياني؟

### Charts: Data Prompt

""You are an expert in chart data extraction. You are given a chart image and you should provide the chart data in CSV format. Here are some examples.

Example 1:

```
""csv
النوع الأدبي,المبيعات (بالآلاف)
روايات,350
خيال علمي,120
فنتازيا,180
حياتي,90
تاريخ,70
علم نفس,110
مذكرات,85
تكنولوجيا,160
فنون,45
أطفال,200
""
```

Example 2:

```
""csv
موضوع,نسبة العملاء الإيجابية,نسبة العملاء السلبية
السياسة في الأردن,47,16
الدين والفكر,35,16
العلاقات غير التقليدية,55,65
العنف في القصص,30,70
العربات القوية,50,80
الثقافة الاجتماعية,60,40
التكنولوجيا والمستقبل,65,35
""
```

Not give me the results as in the previous CSV format.""

### Diagrams: Topic Prompt

""You are an expert in detecting chart types. Below are examples of the expected output format:

Example 1:  
treemap

Example 2:  
flowchart

Example 3:  
diagram

Your task is to determine the type of chart shown in the given image.

**\*\*Instructions:\*\***

- **\*\*Respond with only the chart type\*\*** (e.g., 'flowchart', 'sequence').
- **\*\*Do not provide any explanations, descriptions, or additional text.\*\***
- **\*\*Ensure the output strictly follows the format shown in the examples.\*\***

What type of chart is shown in the image?""

Figure 5: Prompts for Different Task Categories.

## Diagrams: Data Prompt

""You are an expert in diagram data extraction. Your task is to analyze the given diagram and generate structured data in JSON format that captures nodes (entities) and edges (relationships).

## Output Format Example:

for flowchart:

```
```json
{
  "nodes": [
    {
      "id": "1",
      "text": "جمع النفايات",
      "description": "جمع النفايات الصلبة من المناطق الحضرية."
    },
    {
      "id": "2",
      "text": "فرز النفايات",
      "description": "فرز النفايات إلى مواد قابلة لإعادة التدوير وغير قابلة."
    },
    {
      "id": "3",
      "text": "نقل النفايات",
      "description": "نقل النفايات غير القابلة للتدوير إلى مرافق التحويل."
    }
  ],
  "edges": [
    {
      "from": "1",
      "to": "2",
      "text": "فرز"
    },
    {
      "from": "2",
      "to": "3",
      "text": "نقل"
    },
    {
      "from": "3",
      "to": "4",
      "text": "معالجة"
    }
  ]
}
```
```

treemap:

```
```json
{
  "تخصيص التمويل الحكومي والمساعدات": {
    "التنمية الاقتصادية": {
      "دعم المشاريع الصغيرة",
      "تمويل الأثر المناهضة",
      "تحفيز الاستثمار المحلية",
      "إعفاءات ضريبية"
    },
    "البنية التحتية": {
      "تحسين النقل العام",
      "تحديث المحطات والقطارات",
      "الصرف الصحي والمياه",
      "معالجة مياه الصرف"
    }
  }
}
```
```

class diagram:

```
```json
{
  "ورقة_عمل": {
    "خصائص": [
      "اسم الورقة",
      "تاريخ الورقة",
      "مدة الورقة"
    ],
    "علاقات": {
      "تحتوي_على": "جلسة_تدريب",
      "ينظم_من": "مخطط",
      "يشارك_في": "مشارك"
    }
  },
  "جلسة_تدريب": {
    "خصائص": [
      "اسم الجلسة",
      "مدة الجلسة",
      "محتوى"
    ],
    "علاقات": {
      "يقدم_من": "مدرب",
      "ينبع_إلى": "ورقة_عمل"
    }
  }
}
```
```

## Table: HTML Prompt

""Extract the data from the table below and provide the output in HTML format. Output only the data as HTML and nothing else. Here is one example:

```
```html
<table>
<thead>
<tr>
<th>الفئة</th>
<th>النسبة المئوية</th>
<th>التفاصيل</th>
</tr>
</thead>
<tbody>
<tr>
<td>الأهم المحلية</td>
<td>50%</td>
<td>شركة سايك، شركة الاتصالات السعودية، شركة أرامكو</td>
</tr>
<tr>
<td>الأوراق المالية الحكومية</td>
<td>20%</td>
<td>حكومة السعودية، حكومة الإمارات</td>
</tr>
<tr>
<td>السندات الدولية</td>
<td>10%</td>
<td>بنك سويسري، بنك جي بي مورغان</td>
</tr>
<tr>
<td>العقارات التجارية</td>
<td>15%</td>
<td>دي، الرياض، المنامة</td>
</tr>
<tr>
<td>الاستثمارات البديلة</td>
<td>10%</td>
<td>صناديق الاستثمار الخاصة، صناديق التحوط</td>
</tr>
<tr>
<td>التقن وما يعادله</td>
<td>0%</td>
<td>بنك الإمارات دبي الوطني، بنك أبوظبي الأول</td>
</tr>
</tbody>
</table>
```
```

Now generate the data for the provided table.

## Table: Dataframe Prompt

""Extract the data from the table below and provide the output in CSV format. Output only the data as CSV and nothing else. Here is one example:

```
```csv
اسم الشركة,الصفحة,مبلغ الصفقة (مليون دولار),تاريخ الاتفاقية,نوع التقنية
أوراكل,الاستحواد على شركة سيريزو,15-06-28,2023,الحوسبة السحابية
والمنذجة الحيوية
أمازون ويب سيرفيسز,شراكة مع شركة مودلينغ
بيو,20-04-15,2023,المنذجة الحيوية
مايكروسوفت,شراكة مع شركة بيومادكس,10-03-12,2023,الحوسبة
السحابية
جوجل كلاود,شراء شركة بيوكيم سوليوشنز,01-09-35,2023,المنذجة
الحيوية
آي بي إم,توسع في شراكها مع شركة جينوميك
سوفتوير,05-05-18,2023,حوسبة بيولوجية
```
```

Now generate the data for the provided table.

Figure 6: Prompts for Different Task Categories (Continued).