

OCR Hinders RAG: Evaluating the Cascading Impact of OCR on Retrieval-Augmented Generation

Junyuan Zhang^{1*} Qintong Zhang^{1*} Bin Wang^{1*} Linke Ouyang¹ Zichen Wen^{1,4}

Ying Li⁵ Ka-Ho Chow³ Conghui He^{1†} Wentao Zhang²

¹Shanghai AI Laboratory ²Peking University ³The University of HongKong

⁴Shanghai Jiaotong University ⁵Beihang University

Abstract

Retrieval-augmented Generation (RAG) enhances Large Language Models (LLMs) by integrating external knowledge to reduce hallucinations and incorporate up-to-date information without retraining. As an essential part of RAG, external knowledge bases are commonly built by extracting structured data from unstructured PDF documents using Optical Character Recognition (OCR). However, given the imperfect prediction of OCR and the inherent non-uniform representation of structured data, knowledge bases inevitably contain various OCR noises. In this paper, we introduce OHRBench, the first benchmark for understanding the cascading impact of OCR on RAG systems. OHRBench includes 8,561 carefully selected unstructured document images from seven real-world RAG application domains, along with 8,498 Q&A pairs derived from multimodal elements in documents, challenging existing OCR solutions used for RAG. To better understand OCR’s impact on RAG systems, we identify two primary types of OCR noise: Semantic Noise and Formatting Noise and apply perturbation to generate a set of structured data with varying degrees of each OCR noise. Using OHRBench, we first conduct a comprehensive evaluation of current OCR solutions and reveal that none is competent for constructing high-quality knowledge bases for RAG systems. We then systematically evaluate the impact of these two noise types and demonstrate the trend relationship between the degree of OCR noise and RAG performance. Our OHRBench, including PDF documents, Q&As, and the ground truth structured data will be released to foster the development of OCR tailored to RAG and RAG systems that are resilient to OCR noise.

1. Introduction

Retrieval Augmented Generation (RAG) enhances Large Language Models (LLMs) by integrating external knowledge [15, 23], enabling them to respond accurately to queries beyond their training corpus, such as recent news or proprietary content, and reducing hallucinations [21, 23, 34]. This is achieved through a retrieval-then-grounding approach, where relevant documents are retrieved from external knowledge bases and incorporated into the LLM’s prompt for grounding.

As an essential component of RAG systems, the knowledge base defines the scope and quality of documents that RAG can access. Given that a vast amount of real-world knowledge resides in unstructured documents, such as scanned PDFs, constructing an external knowledge base often relies on Optical Character Recognition (OCR) ¹ to parse structured data from these unstructured PDF documents [18, 48]. For instance, MinerU [41] takes raw PDFs as input and extracts plain text, formulas, and tables into structured formats for subsequent RAG applications. However, imperfect predictions of OCR and non-uniform representations of parsing results impair the construction of a high-quality knowledge base for RAG. To be specific, despite advancements in OCR [3, 41, 43], even the leading model cannot achieve perfect accuracy across all scenarios [28, 48]. Furthermore, structural data like table can inherently be parsed in different representation, such as Markdown or LaTeX. These issues introduce OCR noise in parsing results and diminish the quality of the knowledge base. Considering RAG is sensitive to input noise [9, 13, 45], recent works race on downstream RAG components, including more precise retrievers [5, 24, 30] and more advanced LLMs [1, 11, 13, 46]. However, the quality of OCR-based external knowledge bases and its cascading impact on these downstream RAG components have received less attention, which highlights

* These authors contributed equally to this work.

† Corresponding author (heconghui@pjlab.org.cn).

¹We employ the General OCR concept for document parsing from GOT-OCR2.0 [43], which includes, text recognition, multimodal data extraction (table, formula, and chart recognition), and reading order restoration.

a critical but unaddressed gap: *the absence of benchmarks to assess OCR’s cascading impact on each component and entire system of RAG.*

Existing benchmarks either evaluate RAG holistically without fine-grained assessment [49], consider limited OCR solutions without accounting for the noise they introduce [14, 18]. Additionally, they lack documents that present more diverse OCR challenges, such as scanned historical, multilingual, and handwritten documents. To fill this gap, we introduce OHRBench, a question-answering benchmark designed to evaluate OCR’s cascading impact on each component and entire systems of RAG in two ways. First, we construct a document-based RAG Q&A dataset comprising complex, unstructured PDF documents from 7 RAG real-world application areas: *Textbook, Law, Finance, Newspaper, Manual, Academic and Administration*. As detailed in Tab. 1 and Fig. 1, we have collected 8,561 document images featuring attributes that challenges the creation of high-quality knowledge bases for RAG systems. We also provide diverse Q&A pairs which not only span realistic RAG tasks, including understanding, reasoning, and multi-page questions, but also features evidence sourced from key components of OCR in document parsing, making them ideal for assessing the OCR’s impact on RAG performance. Second, we identify two primary OCR noise types: *Semantic Noise*, resulting from prediction errors, and *Formatting Noise*, arising from diverse document element representation. By systematically introducing these noise types into documents, we generate perturbed structured data with varying degrees of noise, enabling further exploration of the quantitative relationship between OCR noise and RAG performance.

With OHRBench, we first conduct a comprehensive benchmark on current OCR solutions, including pipeline-based OCR systems [32, 41], end-to-end OCR models [3, 43] and Vision-Language Models (VLMs) for OCR [2, 6, 7, 42]. We reveal that even the best OCR solutions exhibit a performance gap of 14% at least, compared to the ground truth structured data, facilitating the importance of mitigating OCR noise in RAG systems. Further experiments on different types of OCR noise uncover that *Semantic Noise* consistently exert a significant impact, while *Formatting Noise* affects specific retrievers and LLMs differently, offering valuable insights for developing RAG-tailored OCR solutions and noise robust models.

Contributions. We summarize our main contributions:

- We present OHRBench, a question-answering benchmark designed to evaluate the impact of OCR on RAG systems. OHRBench includes various unstructured PDF documents from seven RAG domains with ground truth structure data annotations and Q&A pairs spanning multiple RAG tasks with diverse source of evidences, posing challenges to the employment of current OCR solutions in RAG systems.
- We conduct a comprehensive evaluation of current OCR

solutions and reveal that none of them is competent for constructing high-quality knowledge bases for RAG systems.

- We identify two primary types of OCR noise, including *Semantic Noise* and *Formatting Noise*, generate perturbed data with varying levels of noise and explore the trend relationship between the degree of OCR noise and RAG performance.

2. Related Works

2.1. Retrieval-Augmented Generation

Retrieval-Augmented Generation (RAG) [21, 23, 34] integrates external knowledge into large language models (LLMs) to mitigate hallucinations. Although RAG technology enhances the generation capabilities of LLMs, it is notably sensitive to input noise. InfoRAG [45] characterizes this noise in RAG as incorrect and irrelevant content within retrieved text and reveals its impact on RAG performance. RAAT [13] further expands noise into relevant noise, counterfactual noise, and irrelevant types. However, these studies focus solely on chunk-level noise introduced during the retrieval stage and its effect on the generation capabilities of LLMs, leaving the impact of noise derived from OCR results unexplored. GARAG [9] examines typographical errors, a form of OCR noise, but its scope is limited to plain text using only synthetic data, overlooking the variety of OCR noise encountered in real-world RAG applications. In this paper, we reveal the impact of noise introduced during the OCR stage, offering a comprehensive analysis of its impact on RAG systems.

2.2. Document parsing with OCR

OCR-based document parsing is a promising solution for structured data extraction from unstructured documents, facilitating applications like RAG. Current OCR solutions can be summarized into three categories, pipeline-based systems [32, 41], end-to-end models [3, 27, 43], and employing VLMs for OCR [7, 19, 42]. Pipeline-based systems decompose OCR into multiple subtasks, such as layout detection, text, formula, and table recognition, enabling fine-grained data extraction. End-to-end models take document images as input and output the overall recognition result in an end-to-end manner. Due to the achievement of VLMs on visual understanding, recent works have explored its application in OCR [28]. In this paper, we evaluate these OCR paradigms, examining their suitability for RAG applications across diverse, real-world document domains.

2.3. Benchmark and evaluation of Retrieval-Augmented Generation

Frameworks like RAGAS [12] and ARES [36] propose evaluating RAG systems based on context relevance, answer

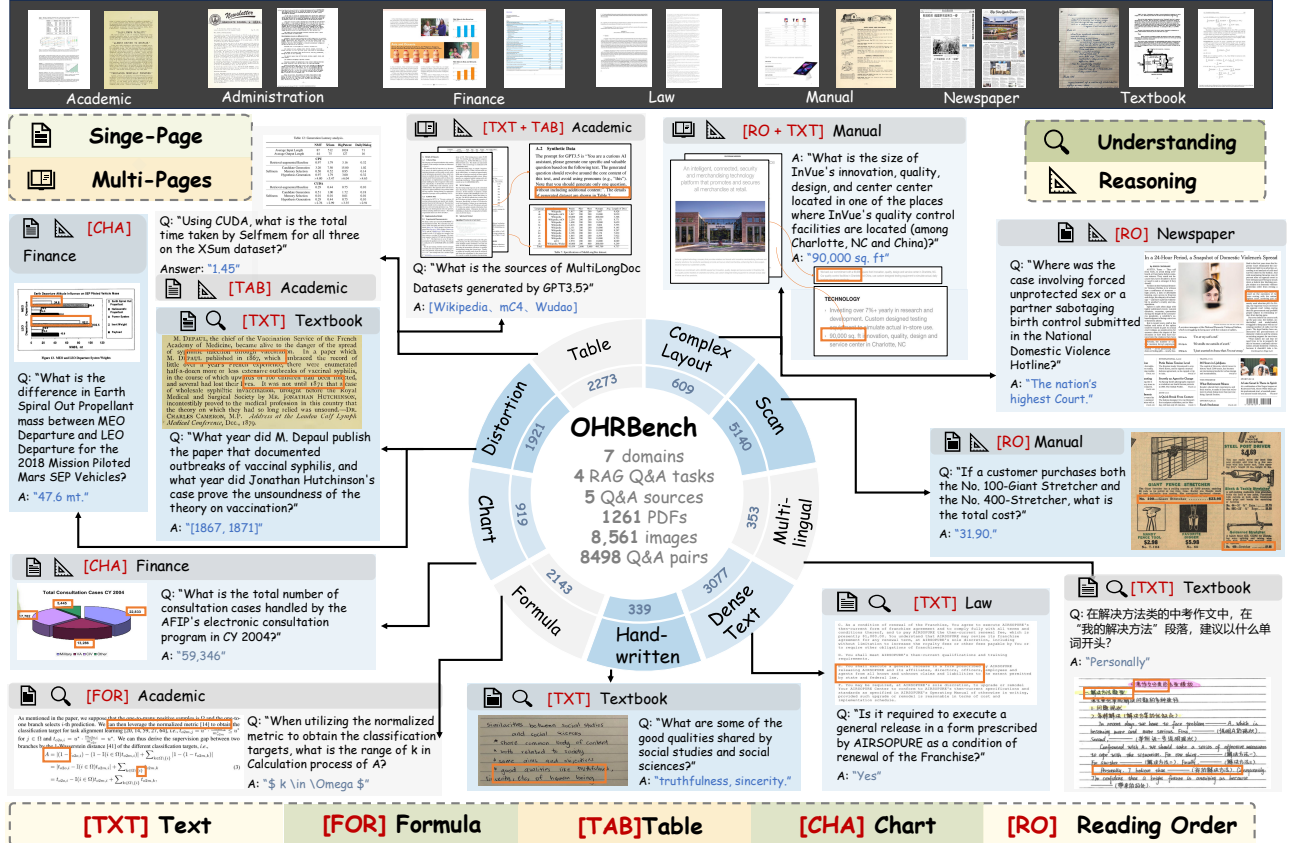


Figure 1. Our OHRBench comprises documents from 7 domains, 9 challenging attributes for OCR, 4 types of Q&A tasks, and 5 Q&A evidence sources. Each number indicates the count of PDF pages with that attribute. Criteria for these attributes can be found in Appendix Sec. II.3

faithfulness, and answer relevance, using LLMs or fine-tuned discriminators for measurement. RGB [4] assesses the noise robustness, negative rejection, information integration, and counterfactual robustness of RAG in news data. MultiHop-RAG [37] focuses on multi-hop reasoning capabilities, while ClashEval [44] explores the context preference in conflicting evidence scenarios. However, these evaluations target specific components of RAG systems, and none of them discusses the impact of external knowledge base construction on RAG systems. Although UDA [18], VisRAG [47] and M3DocRAG [8] explore RAG’s effectiveness in document understanding, they consider limited OCR solutions, ignoring challenging documents and lacking analysis of different OCR noise types. In this paper, we introduce OHRBench to comprehensively investigate OCR noise’s impact on RAG systems.

3. OHRBench

Our OHRBench consists of 1) a number of unstructured PDF documents from seven real-world RAG applications, Q&A pairs derived from multimodal document elements and ground truth structure data annotation, and 2) perturbed structured data based on ground truth with varying de-

grees of OCR noises. Fig. 2 illustrates the construction of OHRBench. We will now delve into the details of each component.

3.1. Data collection

According to [25, 48], extracting structured data from multimodal document elements like formulas and tables poses significant challenges to current OCR solutions. Considering the practical application scenarios of RAG and the challenging field of OCR, we compile a PDF document collection representing seven common RAG application scenarios: *Textbook*, *Law*, *Finance*, *Newspaper*, *Manual*, *Academic* and *Administration*. This collection includes a diverse array of documents from both existing datasets and public web resources. Specifically, we first collect PDF documents from a wide range of existing datasets, including DUDE [40], OmniDocBench [31], FinanceBench [20], CUAD [16], and GNHK [22]. This results in a highly diverse PDF dataset that encompasses complex structured data and layouts, high text density, handwritten, scanned, and historical documents, as well as multilingual content (Chinese and English), which covers most challenges faced by OCR in document parsing. In addition to existing datasets,

Statistic	Number
Documents	1,261
- Domain	7
- Total Pages	8,561
- Avg.Tokens	1,034/page
- Avg.Data Type	1.9/page
Questions	8,498
Avg.Question Token	17.9 ± 8.7
Avg.Answer Token	3.5 ± 4.1
(Evidence Source)	
- Text	3,528 (41.5%)
- Table	2,364 (27.8%)
- Formula	1,267 (14.9%)
- Chart	768 (9.0%)
- Reading Order	691 (8.1%)
(Answer Format)	
- String	4,171 (54.4%)
- Numeric	3,004 (35.3%)
- Yes/No	594 (7.0%)
- List	483 (6.9%)
(Task Type)	
- Understanding	6,114 (71.9%)
- Reasoning	2,384 (28.1%)
- Single-page	7,656 (90.1%)
- Multi-page	842 (9.9%)

Table 1. Dataset Statistics

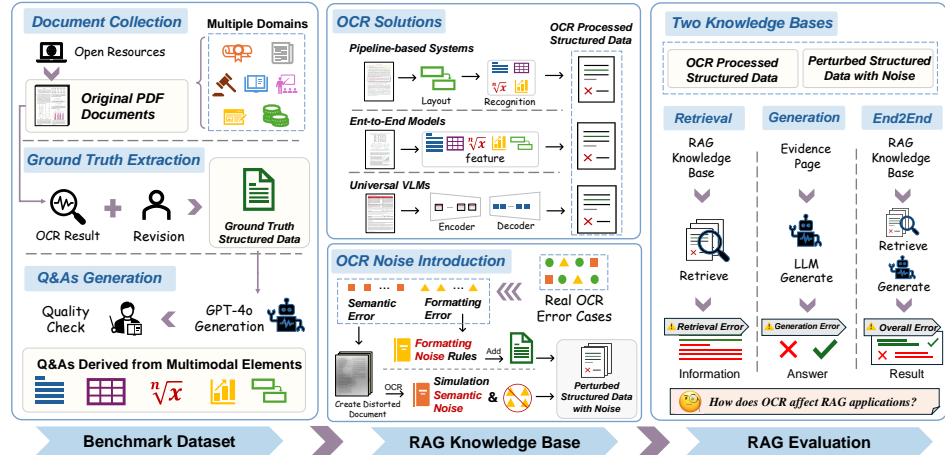


Figure 2. Construction of OHRBench and evaluation protocol. (1) Benchmark Dataset: documents from seven domains, human-verified ground truth structured data, and Q&As from multimodal document elements. (2) RAG Knowledge Base: Current OCR results for benchmarking and perturbed data for assessment. (3) Evaluation of OCR impact on each component and the overall RAG system.

we further supplement our collection with documents from public resources to balance the distribution. We filter out the corrupted or license-restricted documents and finally curate a document dataset comprising 1,261 PDFs and 8,561 images. For each collected document, we manually categorize them into 7 domains and provide ground truth structured data. Specifically, we begin with parsing all documents using Mathpix² for structural data extraction. We then ask expert-level annotators to revise the results, ensuring fidelity to the original structure and content of PDFs while mitigating any style deviations from Mapthix. Detailed descriptions of our selection and processing pipeline can be found in the Appendix Sec. II.

3.2. Q&A pairs generation

The process of extracting structured data from documents involves three key tasks: recognizing plain text; extracting multimodal document elements, including tables, formulas, and charts; and restoring reading order which includes multi-column and truncated paragraph merging. To systematically assess the impact of OCR results on RAG performance, our Q&A generation approach revolves around these 5 evidence sources and various realistic Q&A tasks. Specifically, we provide the ground truth structured data of each document page to GPT-4o and prompt it to generate Q&A based on important components in document parsing, including plain text, tables, formulas, and charts. For questions related to reading order, we identify paragraphs that require merging and instruct GPT-4o to create questions that necessitate combining these paragraphs for a complete answer. We generate both understanding questions, which only require extracting specific information, and reasoning questions, which involve

arithmetic operations, comparisons, or synthesizing information across multiple sections. For multi-page Q&A, we derive them from both single-page Q&A and ground truth structured data from different pages that share the same entity name, recognized with spaCy [17]. Detailed process and the prompt template for the Q&A generation are provided in Appendix Sec. I. Each Q&A consists of the following fields: one page of the original PDF document, evidence context from ground truth structured data that provides the answer to the question, type of evidence (plain text, table, formula, chart and reading order), and the question and answer which are both derived from this evidence context. In this way, these Q&As can serve as a testbed for evaluating OCR results on multimodal document elements.

Quality Control. The quality of Q&A pairs generated by a large language model (LLM) can vary significantly. To address this issue, we apply three data selection criteria to ensure high-quality Q&As: (1) compatibility with realistic RAG applications, (2) faithfulness to task definitions and (3) correctness. We incorporate both heuristic methods and prompting LLM for auto data filtering:

- **Compatibility to RAG Applications.** Questions should be context-independent and not answerable by the model’s internal knowledge. We collect keywords from existing context-dependent questions, such as “according to the document”, as heuristic rules. Following [47], we also employ LLMs to classify context dependence for further filtering. Questions answerable without retrieval are excluded by instructing LLMs to answer without access to the evidence context.
- **Faithfulness to Task Definition.** We ensure questions align with their task definition (reasoning or understanding) using LLMs to judge and that evidence sources match

²<https://mathpix.com/>

the context using heuristic rules. In addition, for multi-page question, we provide single evidence context and use LLMs’ responses to filter out answerable questions.

- **Correctness.** We verify the accuracy of both evidence context and answers. In specific, we provide LLMs with oracle evidence contexts and sample answers repeatedly to filter Q&A below a certain correctness threshold.

Finally, we manually check against the criteria to ensure the quality. The LLMs used in our verification include GPT-4o and DeepSeek-V3 [26], where we find that DeepSeek-V3 achieve similar performance compared to GPT-4o on Q&A verification. A detailed description can be found in the Appendix Sec. II.5. This multi-step quality control ensures the Q&A dataset meets diverse evidence source requirements and practical RAG applications. Ultimately, we filter out 8498 high-quality Q&As from 15317 candidates.

3.3. Data perturbation with OCR noise

Despite advancements in OCR, real-world applications often encounter document types beyond the training corpus of OCR models, leading to low-quality data extraction. Additionally, the different structured representations of document elements further introduces noise, impacting RAG performance. In this paper, we focus on two key types of OCR noises: *Semantic Noise* and *Formatting Noise*. To quantitatively analyze their effects, we start from errors in current OCR results and generate perturbed data with different noise levels. We will now delve into the details of each type.

Semantic Noise results from OCR prediction errors that impact the semantics of parsed content, deviating retrievers and LLMs from integrating correct information related to user queries. To systematically capture realistic *Semantic Noise*, we include diverse perturbation to document images and utilize multiple OCR solutions to perform OCR on these document images, capturing a wide range of real-world *Semantic Noise* as much as possible. We begin with collecting naturally distorted documents and identifying common degradation patterns, such as background artifacts, watermarks, and structural distortions (e.g., dilation and erosion). Then, we extend from [3], where its method has been shown to be effective for simulating naturally distorted documents, we refine perturbation strategies through an iterative, cross-validated process involving multiple annotators. One annotator adjusts distortion parameters and applies them to document images, while a another annotator, who is unaware of the applied modifications, distinguish which document appears artificially altered. This refinement continues until the perturbations become indistinguishable from real samples. Through this process, we identify 8 effective perturbation types that balance realism and intensity. Details and examples are provided in Appendix Sec. III.2. By varying the number and type of perturbations, we generate 3 distinct datasets with controlled *Semantic Noise* levels. We

then choose MinerU, GOT, and Qwen2.5-VL to curate 9 perturbed data with diverse appearance of *Semantic Noise*, enabling a systematic evaluation of its impact on RAG performance.

Formatting Noise stems from stylistic commands, such as white space characters for beautifying formulas and bold and italic commands for better readability, and inconsistencies in structured data representations across Markdown, LaTeX, and HTML. Although irrelevant to semantics, this noise complicates information integration for both retrievers and LLMs. To assess the impact of *Formatting Noise* on RAG, we identify common OCR-induced formatting inconsistencies and develop heuristic rules to introduce controlled perturbations through additions, removals, and format conversions. A detailed list of perturbation rules is in Appendix Sec. III.1. By applying these modifications at varying proportions in ground truth structured data, we create three datasets with different degrees of *Formatting Noise*. Additionally, we evaluate RAG performance under different structured data formats, comparing retrieval and reasoning consistency across Markdown, LaTeX, and HTML representations.

4. Experiments

4.1. Experimental settings

We evaluate the impact of OCR on RAG systems in three ways: retrieval performance, generation performance, and overall system performance. For the retrieval stage, we utilize knowledge bases derived from the same domains of user queries and retrieve the top-2 matched chunk. During the generation stage, we provide the page where the question is derived from for LLMs to generate the response. In the overall evaluation, retrievers retrieve the relevant chunks from the knowledge base in the same domain as the question, and LLMs generate responses based on these chunks. In the overall evaluation, we provide the top-2 matched chunk for generation unless otherwise stated. The default chunk size is 1024 with no overlap.

Metrics. To evaluate the quality of OCR results, we calculate the edit distance between each page of OCR results and the ground truth structured data and report the average values. For assessing retrieval performance, as results of different OCRs often include various extraneous characters, discriminating whether the evidence exactly appears in the retrieved contents is not fair. Following [18], we employ Longest Common Subsequence (LCS) [33] to measure evidence inclusion in retrieved content. For the generation stage, we employ the F1-score metric to measure the accuracy of LLMs’ responses.

Retrievers. We consider two primary retrievers: (1) BGE-M3 [5], a recent SOTA dense retriever within its size category. (2) BM-25 [35, 39] is a lightweight sparse retriever ranking document based on the query term frequency.

	OCR	Retrieval						Generation						Overall					
	E.D.↓	TXT↑	TAB↑	FOR↑	CHA↑	RO↑	ALL↑	TXT↑	TAB↑	FOR↑	CHA↑	RO↑	ALL↑	TXT↑	TAB↑	FOR↑	CHA↑	RO↑	ALL↑
Ground Truth	-	81.2	69.6	74.8	70.3	9.8	70.0	49.4	46.0	34.0	47.0	28.2	43.9	45.0	34.6	28.0	32.9	18.7	36.1
<i>Pipeline-based OCR</i>																			
MinerU [41]	0.24	67.7	48.5	51.1	16.5	5.9	50.1	45.9	39.3	28.6	9.7	29.5	<u>36.7</u>	41.4	28.5	23.0	9.3	17.8	30.0
Marker [32]	0.28	75.2	<u>57.8</u>	<u>55.4</u>	19.7	5.9	56.6	<u>44.5</u>	37.8	27.8	10.9	<u>26.2</u>	35.9	40.1	28.1	22.3	10.0	<u>16.2</u>	29.5
<i>End-to-end OCR</i>																			
GOT [43]	0.27	62.1	41.0	48.7	17.4	3.7	45.4	37.5	28.5	24.1	8.5	7.1	27.8	35.3	22.9	20.1	8.2	5.3	24.6
Nougat [3]	0.34	59.1	32.7	44.2	11.3	4.4	40.9	36.7	22.9	22.9	6.4	6.9	25.5	33.5	18.4	19.4	5.8	3.6	14.5
<i>Vision-Language Model for OCR</i>																			
Qwen2.5-VL-72B [2]	0.18	<u>74.6</u>	59.8	59.7	<u>38.2</u>	5.3	59.2	44.4	42.1	31.8	27.0	11.6	37.5	<u>40.6</u>	31.1	26.1	<u>19.0</u>	8.8	31.1
InternVL2.5-78B [6]	0.28	68.2	57.7	55.3	45.1	2.7	55.8	41.8	<u>41.8</u>	<u>29.0</u>	33.6	3.3	35.8	38.2	<u>31.0</u>	<u>23.3</u>	22.9	3.1	29.6

Table 2. Evaluation of various OCR solutions and their impacts on RAG systems. The OCR performance is reported using edit distance (E.D.). We report the generalized LCS or F1 of five types of evidence sources, including plain text (TXT), table (TAB), formula (FOR), chart (CHA), and reading order (RO). **Bold** indicates the best performance, and underline indicates the second-best performance.

LLMs. We employ three representative open-source LLMs: Qwen2 (Qwen2-7B-Instruct and Qwen2-72B-Instruct) [42] and Llama-3.1 (Llama3.1-8B-Instruct) [11]. A standard prompt template is used to format responses consistently across all LLMs (see Appendix Sec. I). All open-source models are downloaded from Huggingface³, with inference conducted on 8 NVIDIA A100 GPUs.

4.2. Benchmarking current OCR solutions

In this section, we evaluate the suitability of current OCR solutions for real-world RAG applications by conducting comprehensive experiments with our OHRBench. We involve several representative OCR solutions including (1) Pipeline-based OCR, such as MinerU [41] and Marker [32], (2) End-to-end OCR, including GOT [43] and Nougat [3], and (3) Vision-Language Models, specifically Qwen2.5-VL-72B [38] and InternVL2.5-78B [6]. For GOT, we employ its format OCR mode to output structured data. For Qwen2.5-VL-72B and InternVL2.5-78B, we prompt them to produce formulas, tables, and charts in LaTeX format, with the prompt template available in Appendix Sec. I. The retrievers used are BGE-M3 and BM25, while the LLMs are Llama-3.1-8B-Instruct and Qwen2-7B-Instruct. All metrics are averaged across domains and combinations of retrievers and LLMs. Details of experimental results are available in Appendix Sec. IV.1.

Through the comparison presented in Tab. 2, we derive several key conclusions about the performance of these OCR solutions and their corresponding impacts on RAG systems, as follows: **(1) VLMs for OCR achieve the best overall performance.** Among all OCR solutions, Qwen2.5-VL-72B consistently outperforms others across all three evaluation stages. Its superiority stems from its ability to handle structured data more effectively than both pipeline-based and end-to-end OCR methods. Despite claims that GOT can parse charts, its performance remains subpar. Similarly, MinerU, Marker, and Nougat fail to produce comparable results due to their inability in parsing chart. For plain text

questions, although its poor performance on high-resolution documents (see newspaper in APPendix Sec. IV.1, Qwen2.5-VL-72B performs comparably to pipeline-based OCR. Our manual review suggests that its strong language decoder enhances robustness against historical and distorted documents, a capability lacking in pipeline-based OCR. **(2) Reading order is challenging for VLMs and End-to-end OCR.** Despite their strong semantic understanding capabilities offered by language decoders, both VLM and end-to-end OCR struggle with merging paragraphs correctly, reflected in F1-score of just 8.8 in overall evaluation. In contrast, pipeline-based OCR, though lacking semantic understanding, achieves performance close to ground truth (GT) using rule-based strategies. However, GT itself performs poorly, likely because reading-order questions require integrating information across multiple paragraphs, posing challenges for current RAG systems [37]. **(3) All OCR solutions exhibit performance degradation.** Even the best solutions experience a 14% (5 F1-score) drop in the overall stage evaluation, with greater losses in the retrieval and generation stages. This indicates that our OHRBench presents significant challenges for both OCR solutions and RAG systems.

In summary, current OCR solutions struggle to maintain robustness and effectiveness across diverse real-world RAG application scenarios. Additionally, standard OCR metrics like edit distance do not always align with RAG performance. For example, while MinerU and Qwen2.5-VL-72B exhibit lower edit distances compared to Marker, they do not consistently achieve better performance across all metrics. This discrepancy may be attributed to the varying types of OCR noise introduced by different solutions. To further investigate, we systematically explore the impact of these OCR noise types on RAG in Sec. 4.3.

4.3. In-depth analysis of OCR noise’s impact on RAG

In this section, we conduct an in-depth analysis of the impact of *Semantic Noise* and *Formatting Noise* on RAG systems, using perturbed structured data with varying levels of perturbations. For each type of OCR noise, we introduce three

³<https://huggingface.co/>

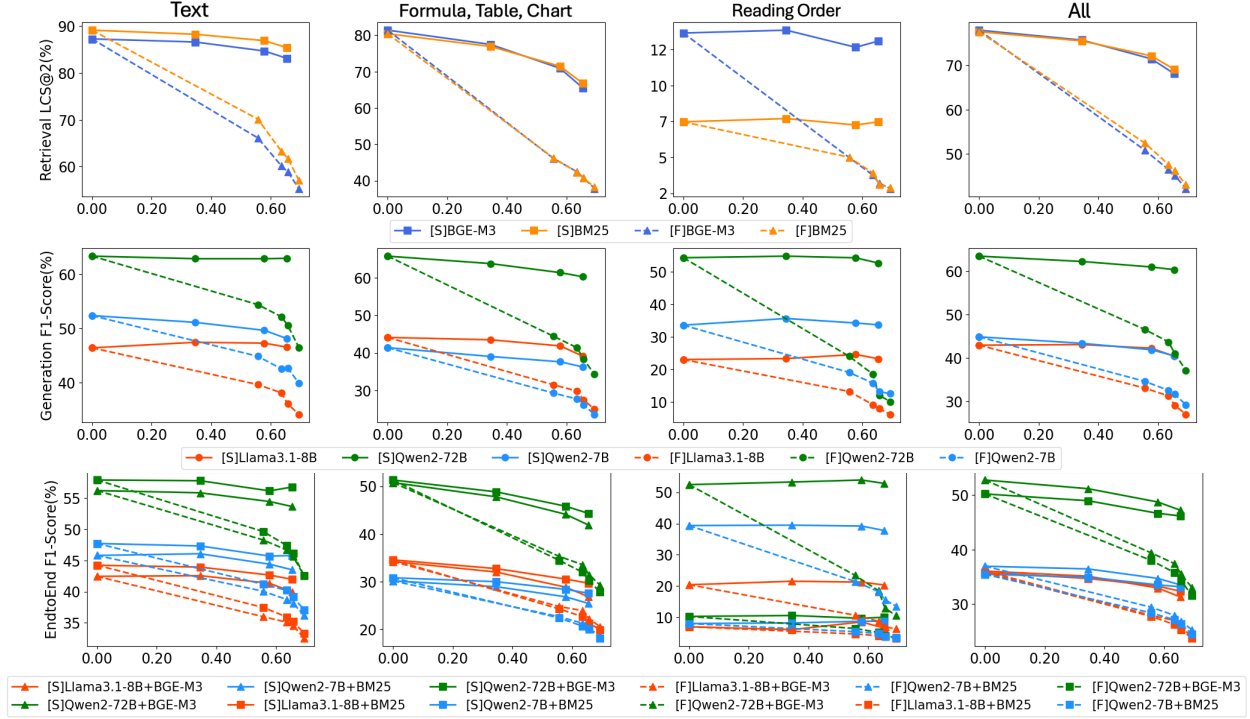


Figure 3. Impact of *Semantic Noise* ([S] dashed lines) and *Formatting Noise* ([F] solid lines) on RAG components. The horizontal axis denotes the ratio r_{noise} , where higher values indicate greater OCR-induced noise. We report LCS and F1-score for each evidence source: text (first column), the average score for multimodal elements (tables, formulas, and charts, second column), reading order (third column), and all sources combined (last column).

noise levels—mild, moderate, and severe—to systematically assess their effects. Since our findings in Sec. 4.2 indicate that edit distance fails to accurately capture the degree of OCR noise, we instead define the ratio r_{noise} , representing the proportion of Q&A pairs affected by OCR noise as a measure. Specifically, for each Q&A pair, we compute the LCS between its evidence context and the corresponding perturbed structured data. If LCS exceeds 0.95, the Q&A pair is considered as unaffected, otherwise, it is affected. We then use the ratio (r_{noise}) to quantify the degree of perturbations, with 0 representing to ground truth structured data and values approaching 1 indicating greater perturbation. This approach allows us to align the perturbation levels of *Formatting Noise* and *Semantic noise* for fair comparisons. Additionally, for *Formatting Noise*, we evaluate its retrieval performance with modified LCS calculations by excluding stylistic commanding introduced during perturbation, ensuring a fair assessment of retrieval accuracy. For each degree of *Semantic Noise*, we report the average RAG performance using three different OCR results, including MinerU, GOT, and Qwen2.5-VL-72B.

4.3.1. Fine-grained impact on retrieval and generation

***Semantic Noise* significantly influences both retrieval and generation phases.** As illustrated in Fig. 3, increasing *Semantic Noise* from mild ($r_{\text{noise}} = 0$) to severe ($r_{\text{noise}} > 0.6$) results in nearly a 50% performance decline for most retriev-

ers and LLMs. In the retrieval stage, both the sparse retriever BM25 and the dense retriever BGE-M3 suffer consistent performance declines across all types of questions, suggesting that dense retrieval’s stronger comprehension does not provide robustness against *Semantic Noise*. In the generation phase, all LLMs struggle with *Semantic Noise*, among which performance on reading-order decreases the most. Interestingly, although the way we introduce *Semantic Noise* should primarily affect text recognition, questions related to multimodal elements (tables, formulas and charts) degrade even further, highlighting the challenges in parsing, understanding and reasoning over multimodal document data.

***Formatting Noise* primarily affects multimodal questions.**

While performance on plain text queries and reading-order-related questions remain largely unaffected, retrieval and generation performance drops more severe for multimodal queries. The maximum performance losses reach 12.7% for BGE-M3 and 9.1% for Llama3.1-8B in retrieval and generation, respectively. In addition, larger LLMs exhibit greater robustness, with only a 7% performance reduction on multimodal questions, indicating that more advanced LLMs can effectively handle *Formatting Noise*.

4.3.2. Impact on end-to-end evaluation

***Semantic Noise* consistently demonstrate a strong impact, while *Formatting Noise* affects specific retrievers and LLMs differently.** *Semantic Noise* consistently de-

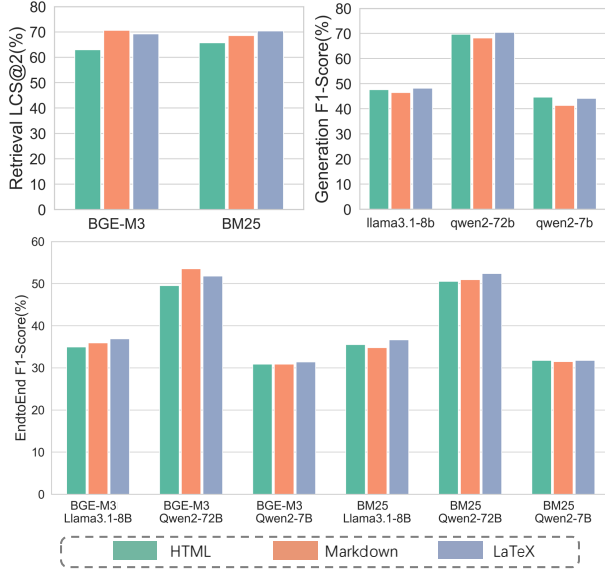


Figure 4. Performance of retrieval, generation and end-to-end with different table format. We only report the results of table-related questions.

grades performance across all combination of retrievers and LLMs, particularly on multimodal questions involving tables, formulas, and charts. In contrast, the effect of *Formatting Noise* is more variable when using smaller LLMs, Llama3.1-8B and Qwen2-7B, for generation, despite greatly reduced retrieval accuracy, the overall performance shows a slight change due to their limited information integration capabilities. Conversely, using a larger LLM, Qwen2-72B, is highly sensitive to retrieval performance. While its generation performance shows light changes, the overall performance decreases more, especially on questions related to multimodal elements.

In summary, *Semantic Noise* significantly affects each stage of RAG and the entire system. The impact of *Formatting Noise* varies with different retrievers and language models, particularly affecting questions related to multimodal elements.

4.3.3. Impact of table format

In addition to perturbations, we investigate the influence of different table formats as a kind of *Formatting Noise*. As shown in Fig. 4, HTML tables show inferior performance during retrieval compared to the Markdown and LaTeX formats. Markdown and LaTeX formats perform similarly, with BGE-M3 demonstrating a better understanding of Markdown. In the generation phase, HTML and LaTeX showed similar performance across all models, but the Markdown format performed worse due to the lack of support for merging cells. In end-to-end evaluations affected by low retrieval performance, using HTML tables is comparatively worse, while the combination of Qwen2-72B and BGE-M3 using

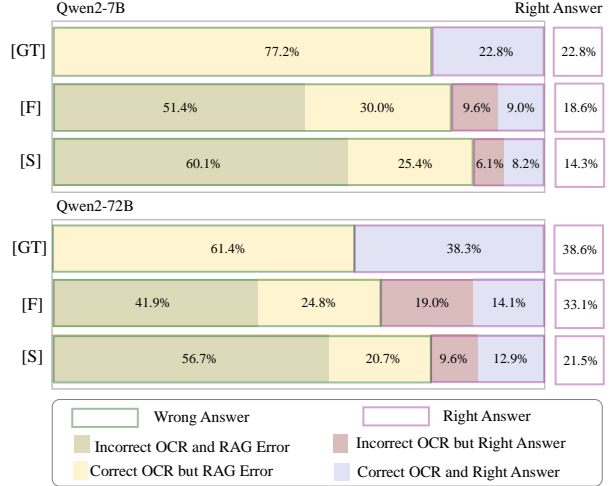


Figure 5. Analysis of answer correctness distribution in Q&A pairs, using different perturbed data and LLMs. It reveals that larger LLMs are more robust to OCR noise. [S] and [F] denote perturbed data with severe *Semantic Noise* and *Formatting Noise*, respectively. [GT] represents ground truth perturbed data.

Markdown achieves the best performance.

4.3.4. Error analysis

We further conduct error analysis to understand the bottlenecks of the RAG system with OCR noise in a quality approach. Specifically, we calculate the distribution of OCR errors and RAG errors when using ground truth structured data and perturbed data with severe *Formatting Noise* and *Semantic Noise*. Our evaluation employs BGE-M3 as the retriever and assesses error distributions using two generators: Qwen2-7B and Qwen2-72B. We use the same strategies to identify the proportion of OCR errors as in Sec. 4.3. The distribution of these errors is illustrated in Fig. 5. It indicates that when the proportion of OCR errors is similar for *Semantic Noise* and *Formatting Noise* (66% vs. 61%), it performs better with *Formatting Noise*. Of these, about half of the correct responses when using Qwen2-7B as a generator were done by the model with incorrect OCR results. Meanwhile, larger LLMs are more robust to OCR noise. Compared to Qwen2-7B, Qwen2-72B has a 9.4% and 3.5% higher percentage of samples with OCR errors but ultimately correct in *Formatting Noise* and *Semantic Noise*, respectively.

5. Conclusion

In this paper, we present OHRBench to evaluate the impact of OCR on RAG systems, which encompasses diverse PDF documents from seven RAG application scenarios along with Q&A pairs derived from multimodal elements in these documents. Through comprehensive evaluations of current OCR solutions, we reveal that none is fully capable of RAG systems across all scenarios. Furthermore, our analysis of different types of OCR noise demonstrates that while no

retrievers and LLMs are immune to *Semantic Noise*, more advanced models exhibit greater resilience to *Formatting Noise*. We believe that our documents featuring challenging OCR attributes and Q&A pairs sourced from varied document elements, will advance the development of OCR solutions tailored for RAG and OCR noise-resistant RAG systems.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1
- [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 2, 6
- [3] Lukas Blecher, Guillem Cucurull, Thomas Scialom, and Robert Stojnic. Nougat: Neural optical understanding for academic documents. In *The Twelfth International Conference on Learning Representations*, 2024. 1, 2, 5, 6, 3
- [4] Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 17754–17762, 2024. 3
- [5] Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. *arXiv preprint arXiv:2402.03216*, 2024. 1, 5
- [6] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. 2, 6
- [7] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. 2
- [8] Jaemin Cho, Debanjan Mahata, Ozan Irsoy, Yujie He, and Mohit Bansal. M3docrag: Multi-modal retrieval is what you need for multi-page multi-document understanding. *arXiv preprint arXiv:2411.04952*, 2024. 3, 1
- [9] Sukmin Cho, Soyeong Jeong, Jeongyeon Seo, Taeho Hwang, and Jong C Park. Typos that broke the rag’s back: Genetic attack on rag pipeline by simulating documents in the wild via low-level perturbations. *arXiv preprint arXiv:2404.13948*, 2024. 1, 2
- [10] Chao Deng, Jiale Yuan, Pi Bu, Peijie Wang, Zhong-Zhi Li, Jian Xu, Xiao-Hui Li, Yuan Gao, Jun Song, Bo Zheng, et al. Longdocurl: a comprehensive multimodal long document benchmark integrating understanding, reasoning, and locating. *arXiv preprint arXiv:2412.18424*, 2024. 2
- [11] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 1, 6
- [12] Shahul Es, Jithin James, Luis Espinosa-Anke, and Steven Schockaert. Ragas: Automated evaluation of retrieval augmented generation. *arXiv preprint arXiv:2309.15217*, 2023. 2
- [13] Feiteng Fang, Yuelin Bai, Shiwen Ni, Min Yang, Xiaojun Chen, and Ruifeng Xu. Enhancing noise robustness of retrieval-augmented language models with adaptive adversarial training. *arXiv preprint arXiv:2405.20978*, 2024. 1, 2
- [14] Manuel Faysse, Hugues Sibille, Tony Wu, Gautier Viaud, Céline Hudelot, and Pierre Colombo. Colpali: Efficient document retrieval with vision language models. *arXiv preprint arXiv:2407.01449*, 2024. 2, 1
- [15] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2023. 1
- [16] Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. Cuad: An expert-annotated nlp dataset for legal contract review. In *Advances in Neural Information Processing Systems*, 2021. 3, 1
- [17] Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017. 4, 2
- [18] Yulong Hui, Yao Lu, and Huanchen Zhang. Uda: A benchmark suite for retrieval augmented generation in real-world document analysis. *arXiv preprint arXiv:2406.15187*, 2024. 1, 2, 3, 5
- [19] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 2
- [20] Pranab Islam, Anand Kannappan, Douwe Kiela, Rebecca Qian, Nino Scherrer, and Bertie Vidgen. Financebench: A new benchmark for financial question answering. *arXiv preprint arXiv:2311.11944*, 2023. 3, 1
- [21] Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Atlas: Few-shot learning with retrieval augmented language models. *Journal of Machine Learning Research*, 24(251):1–43, 2023. 1, 2
- [22] Alex W. C. Lee, Jonathan Chung, and Marco Lee. Gnhk: A dataset for english handwriting in the wild. In *International Conference of Document Analysis and Recognition (ICDAR)*, 2021. 3, 1
- [23] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020. 1, 2
- [24] Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. Towards general text em-

- beddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*, 2023. 1
- [25] Zichao Li, Aizier Abulaiti, Yaojie Lu, Xuanang Chen, Jia Zheng, Hongyu Lin, Xianpei Han, and Le Sun. Readoc: A unified benchmark for realistic document structured extraction. *arXiv preprint arXiv:2409.05137*, 2024. 3
- [26] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024. 5
- [27] Chenglong Liu, Haoran Wei, Jinyue Chen, Lingyu Kong, Zheng Ge, Zining Zhu, Liang Zhao, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Focus anywhere for fine-grained multi-page document understanding. *arXiv preprint arXiv:2405.14295*, 2024. 2
- [28] Yuliang Liu, Zhang Li, Biao Yang, Chunyuan Li, Xucheng Yin, Chenglin Liu, Lianwen Jin, and Xiang Bai. On the hidden mystery of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023. 1, 2
- [29] Yubo Ma, Yuhang Zang, Liangyu Chen, Meiqi Chen, Yizhu Jiao, Xinze Li, Xinyuan Lu, Ziyu Liu, Yan Ma, Xiaoyi Dong, et al. Mmlongbench-doc: Benchmarking long-context document understanding with visualizations. *arXiv preprint arXiv:2407.01523*, 2024. 1
- [30] Gabriel de Souza P Moreira, Radek Osmulski, Mengyao Xu, Ronay Ak, Benedikt Schifferer, and Even Oldridge. Nv-retriever: Improving text embedding models with effective hard-negative mining. *arXiv preprint arXiv:2407.15831*, 2024. 1
- [31] Linke Ouyang, Yuan Qu, Hongbin Zhou, Jiawei Zhu, Rui Zhang, Qunshu Lin, Bin Wang, Zhiyuan Zhao, Man Jiang, Xiaomeng Zhao, Jin Shi, Fan Wu, Pei Chu, Minghao Liu, Zhenxiang Li, Chao Xu, Bo Zhang, Botian Shi, Zhongying Tu, and Conghui He. Omnidocbench: Benchmarking diverse pdf document parsing with comprehensive annotations, 2024. 3, 1, 2
- [32] Vik Paruchuri. Marker, 2024. 2, 6
- [33] Mike Paterson and Vlado Dančik. Longest common subsequences. In *International Symposium on Mathematical Foundations of Computer Science*, pages 127–142. Springer, 1994. 5
- [34] Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. In-context retrieval-augmented language models. *Transactions of the Association for Computational Linguistics*, 11:1316–1331, 2023. 1, 2
- [35] Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, et al. Okapi at trec-3. *Nist Special Publication Sp*, 109:109, 1995. 5
- [36] Jon Saad-Falcon, Omar Khattab, Christopher Potts, and Matei Zaharia. Ares: An automated evaluation framework for retrieval-augmented generation systems. *arXiv preprint arXiv:2311.09476*, 2023. 2
- [37] Yixuan Tang and Yi Yang. Multihop-rag: Benchmarking retrieval-augmented generation for multi-hop queries. *arXiv preprint arXiv:2401.15391*, 2024. 3, 6
- [38] Qwen Team. Qwen2.5: A party of foundation models, 2024. 6
- [39] Andrew Trotman, Antti Puurula, and Blake Burgess. Improvements to bm25 and language models examined. In *Proceedings of the 19th Australasian Document Computing Symposium*, pages 58–65, 2014. 5
- [40] Jordy Van Landeghem, Rubèn Tito, Łukasz Borchmann, Michał Pietruszka, Paweł Joziak, Rafał Powalski, Dawid Jurkiewicz, Mickaël Coustaty, Bertrand Anckaert, Ernest Valveny, et al. Document understanding dataset and evaluation (dude). In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19528–19540, 2023. 3, 1
- [41] Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, et al. Mineru: An open-source solution for precise document content extraction. *arXiv preprint arXiv:2409.18839*, 2024. 1, 2, 6
- [42] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 2, 6
- [43] Haoran Wei, Chenglong Liu, Jinyue Chen, Jia Wang, Lingyu Kong, Yanming Xu, Zheng Ge, Liang Zhao, Jianjian Sun, Yuang Peng, et al. General ocr theory: Towards ocr-2.0 via a unified end-to-end model. *arXiv preprint arXiv:2409.01704*, 2024. 1, 2, 6
- [44] Kevin Wu, Eric Wu, and James Zou. How faithful are rag models? quantifying the tug-of-war between rag and llms’ internal prior. *arXiv preprint arXiv:2404.10198*, 2024. 3, 1
- [45] Shicheng Xu, Liang Pang, Mo Yu, Fandong Meng, Huawei Shen, Xueqi Cheng, and Jie Zhou. Unsupervised information refinement training of large language models for retrieval-augmented generation. *arXiv preprint arXiv:2402.18150*, 2024. 1, 2
- [46] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024. 1
- [47] Shi Yu, Chaoyue Tang, Bokai Xu, Junbo Cui, Junhao Ran, Yukun Yan, Zhenghao Liu, Shuo Wang, Xu Han, Zhiyuan Liu, et al. Visrag: Vision-based retrieval-augmented generation on multi-modality documents. *arXiv preprint arXiv:2410.10594*, 2024. 3, 4, 2
- [48] Qintong Zhang, Victor Shea-Jay Huang, Bin Wang, Junyuan Zhang, Zhengren Wang, Hao Liang, Shawn Wang, Matthieu Lin, Wentao Zhang, and Conghui He. Document parsing unveiled: Techniques, challenges, and prospects for structured

information extraction. *arXiv preprint arXiv:2410.21169*, 2024. [1](#), [3](#)

- [49] Anni Zou, Wenhao Yu, Hongming Zhang, Kaixin Ma, Deng Cai, Zhuosheng Zhang, Hai Zhao, and Dong Yu. Docbench: A benchmark for evaluating llm-based document reading systems. *arXiv preprint arXiv:2407.10701*, 2024. [2](#), [1](#)

OCR Hinders RAG: Evaluating the Cascading Impact of OCR on Retrieval-Augmented Generation

Supplementary Material

I. Instruction Prompts

Q&A Generation Prompt Template. The template is shown in Tab. S6. Following [44, 49], we instruct GPT-4o to generate questions with clear entities and require three levels of difficulty for question diversity.

RAG Generation Stage Prompt Template. The prompt template for LLMs and VLMs with text-only input is shown in Tab. S9.

Vision-Language Models OCR Prompt Template. We tune the prompt for the best performance of VLMs OCR, by comparing simple and detailed instructions as shown in Tab. S10. Results in Tab. S11 indicate that the detailed prompt consistently performs better across all evaluations, so it is used by default.

II. Benchmark Construction Details

II.1. Document details

We curate a dataset of 1,261 PDFs spanning 8,561 pages, with 3,596 pages designated for Q&A generation and the remainder forming part of the knowledge base. These PDFs are sourced from DUDE [40], OmniDocBench [31], FinanceBench [20], CUAD [16], GNHK [22], and public resources, including Arxiv⁴, ManualsLib⁵, LibreTexts⁶.

DUDE: We extract documents from the validation and test splits of DUDE, applying manual screening based on our criteria Fig. 1 to exclude samples infeasible for structured data parsing and classify each of them into 7 domains. We finally selected 450 PDFs with 4,058 images from 2,069 PDF candidates.

OmniDocBench: OmniDocBench [31] features span-level annotations and presents challenges for OCR due to its multilingual, high-resolution with dense text and handwritten content. We select all newspaper documents and manually review textbook-related samples, eliminating those with low knowledge density that hinder meaningful Q&A generation. This process yields 289 PDFs.

FinanceBench: Following prior observations [29], both DUDE [40] and FinanceBench [20] contain diverse document types. From FinanceBench, we randomly sample 10 PDFs characterized by large, complex tables and charts.

CUAD: We randomly select 65 PDFs to supplement the documents in law domains, which all have high text density.

GNHK: GNHK consists of handwritten documents. We manually assess and remove those with low knowledge density, finalizing a selection of 172 PDFs.

Each document is manually reviewed by primary authors to ensure its availability for academic use. Detailed domain statistics are shown in Tab. S1

Domains	PDFs	Pages	Pages with Q&As
Law	95	1187	1143
Finance	65	2133	1359
Textbook	504	678	1126
Manual	87	1724	1155
Newspaper	279	487	1202
Academic	85	1011	1181
Administration	146	1341	1332
Total	1261	8561	8498

Table S1. Document statistics of each domain

II.2. Ground truth structured data annotation

We annotate the ground truth structured data using Mathpix Markdown format, where tables and formulas are represented in LaTeX. Chart data is extracted in LaTeX table format, with charts lacking clear numeric values in figure filtered out. For images in documents, any parsable text is retained as plain text in the corresponding section. To ensure high-quality annotations, we first use Mathpix to pre-annotate all PDFs. Finally, the primary authors employ Mathpix Markdown previews⁷ to render structured data into PDFs, manually review and correct pre-annotated results.

II.3. Document with challenging attributes

Although existing RAG document benchmarks have gathered PDFs from different domains [8, 14, 18], they often ignore the challenges posed by OCR. To address this, we construct a benchmark that explicitly incorporates documents with challenging attributes. We define nine key attributes: structured data (tables, formulas, charts), complex layouts, handwritten content, distortions, scanned PDFs, dense text, and multilingual content. Structured data, dense text (exceeding 770 tokens), and multilingual pages are classified based on the annotated ground truth structured data. A document is considered to have a complex layout if its layout

⁴<https://arxiv.org>

⁵<https://www.manualslib.com/>

⁶<https://libretexts.org/>

⁷<https://github.com/Mathpix/vscode-mathpix-markdown>

detection yields more than 20 bounding boxes. Distorted, scanned, and handwritten documents are identified during manual checks.

II.4. Q&A generation

To generate high-quality Q&A pairs covering diverse tasks and evidence sources, we define multiple prompts for each task, as detailed Tabs. S6 to S8. For Chinese questions, we provide the same set of templates in Chinese to ensure that the model generates Chinese responses. **Q&A with different evidence sources.** For Q&A generation with evidence sourced from plaintiff text, table, formula and chart, we extract relevant pages from the ground truth structured data and use GPT-4o to generate Q&A pairs grounded in the corresponding evidence via tailored prompts. For Q&A related to reading order, we leverage MinerU [41], the leading model for reading order recognition [31], to identify the reading order and bounding box of paragraphs in each document. When working with documents from OmniDocBench [31], we directly use the ground truth reading order from its annotations. We verify the layout detection and reading order predictions, selecting paragraph pairs that meet one of the following criteria:

- Adjacent paragraphs in reading order whose bounding boxes are not vertically aligned.
- Paragraphs separated by multimodal document elements (e.g., block formulas, tables, or images).

We then randomly sample 1,500 candidate matches, manually correcting approximately 20% where MinerU’s predictions are inaccurate. We then prompt GPT-4o to generate Q&A pairs using the prompts in Tab. S7. We find that this simple prompting-based strategy can effectively generate questions with diverse evidence sources, with over 90% correctly aligned with their evidence source in our Q&A verification process.

Q&A with different tasks. To generate both understanding and reasoning questions, we apply the corresponding prompts from Tab. S7. For multi-page Q&A generation, we employ two different approaches to generate Q&A candidates: (1) Combine questions from two single-page Q&As that mention the same entity. (2) Generating multi-page questions from two paragraphs on different pages that reference the same entity. Specifically, we use spaCy [17] for named entity recognition in both single-page Q&As and document paragraphs. We then filter out candidate pairs, including: (1) Single-page Q&A pairs where the entity in one answer appears in another question. (2) Paragraph pairs that share the same entity. We finally utilize the prompts in Tab. S8 to generate multi-page questions. However, despite the many optimizations of the prompt and generation strategies we tried, GPT-4o sometimes produces Q&A pairs that are either answerable with a single paragraph or simply concatenate two single-page questions while maintaining separate evi-

dence sources instead of high-quality and realistic multi-page Q&As. To address these limitations, we develop a comprehensive filtering process to ensure the quality of multi-page Q&As, as detailed in Sec. II.5.

II.5. Q&A verification.

We verify Q&A quality based on three criteria: (1) Compatibility with realistic RAG applications, (2) faithfulness to task definition, and (3) correctness. Below, we detail our approach for each aspect.

Compatibility with Realistic RAG Applications. To assess context dependence, we identify key patterns from existing context-dependent questions and apply the following heuristics:

- Questions lacking an explicit entity name.
- Questions containing more than one ambiguous pronouns (e.g., "he," "she," "it," "they", "this", "that").
- Questions featuring phrases such as "in the document" or "according to the document."

These rules filter most context-dependent questions. We then refine the selection using prompts in VisRAG [47] and DeepSeek-V3 to further distinguish context-dependent questions from the remaining set. Additionally, we use GPT-4o to exclude questions answerable without retrieval by instructing it to respond without providing evidence context across both single-page and multi-page Q&As.

Faithfulness to Task Definition. Based on the Q&A verification prompts in [10], we use the prompts in Tab. S12 to assess faithfulness using DeepSeek-V3. To verify the validity of evidence sources, we locate them in the original ground truth structured data and ensure they originate from the correct corresponding LaTeX code environments. For the multi-page and reading-order questions, we employ GPT-4o to generate three responses: (1) without context, (2) with context A, and (3) with context B. If any response yields a correct answer, the question is excluded, ensuring that only truly multi-page or reading-order-related questions remain.

Correctness. To guarantee each Q&A has a unique and correct answer supported by its evidence context, we provide oracle evidence and sample GPT-4o’s response 10 times. We apply a best-of-N strategy to determine the final answer, which must match the ground truth. Q&As with fewer than three consistent correct responses are also excluded.

Our filtering pipeline underwent two iterations of refinement. In each round, we randomly sample 100 Q&As to verify the filtering results adherence to our criteria. Finally, to mitigate false positives, we manually reviewed all remaining questions, yielding 8,498 high-quality Q&As from an initial pool of 15,317 candidates.

III. OCR Noise Introduction

III.1. Rules for Formatting Noise introduction

To introduce *Formatting Noise*, we define a perturbation rate r to control its extent. In order to match the level of *Semantic Noise* (measured by similar edit distance), we set the $r = \{0.1, 0.3, 0.6\}$, indicating the three levels of perturbation: mild, moderate, and severe. Based on the *Formatting Noise* in the existing OCR results, we formulate the following rules to perturb plain text, tables, and formulas, respectively.

III.1.1. Plain text

Text Style: Given the plain text content of the ground truth, we randomly divide it into a sequence where each item consists of 2 to 5 words, select target items based on r , and apply one of the following operations as perturbations.

- **Bold:** Enclose the selected text in `**` or `\textbf{}`.
- **Italic:** Enclose the selected text in `*` or `\textit{}`.
- **Underline:** Enclose the selected text in `_` or `\underline{}`.

Title Formatting: We identify short sentences that end with a full stop and have no more than 5 words as potential headings. We randomly pick them according to r and add one of level 1 to level 3 title controls in Markdown (`#`) or LaTeX (`\section{}`) to make new titles.

Paragraph: To simulate the line breaks that exist in PDFs, we randomly insert `\n` at word intervals based on r .

III.1.2. Formula

Formula Conversion: Randomly convert the inline formula into block formula and vice versa at rate r .

Extraneous Elements: We first randomly select the target formulas based on r . Subsequently, for each target formula, we randomly insert 1 to 5 meaningless markers in its symbol gaps, including `\`, `\quad`, `\qquad`, `\;`, `\:`.

Equivalent Symbols: For each formula, we replace the following equivalent symbols with probability r :

- **bold:** `\mathbf{}`, `\boldsymbol{}`.
- **cursive:** `\mathbb{}`, `\pmb{}`, `\mathrsfs{}`, `\euscript{}`, `\mathcal{}`.
- **unicode:** (`\sigma`, `\u03A3`), etc⁸.

III.1.3. Table

Row and Column Lines: For each line and column, randomly insert `\hline` or `\cline` with probability r .

Cell Content: For each cell content, randomly apply above rules for plain text or formula with probability r .

III.2. Rules for Semantic Noise introduction

In order to construct perturbed document images that conform to the realistic distribution of naturally distorted docu-

⁸Full lists are drawn from <https://raw.githubusercontent.com/w3c/xml-entities/refs/heads/gh-pages/unicode.xml>

OCR	Avg. Counts
MinerU	35.0
GOT	45.7
Nougat	63.2
F-minor	37.9
F-moderate	42.2
F-severe	56.3

Table S2. Counts of *Formatting Noise*. The counts of *Formatting Noise* we add (F-minor, F-moderate, F-severe) is approximately the distribution of the counts of *Formatting Noise* for MinerU, GOT and Nougat.

ments, we use a cross-validated process involving multiple annotators. We finally identify 8 strategies from [3] as follows:

- **Background Addition:** We collect 15 background images of real paper textures and blend them with original images at an 80:20 ratio.
- **Salt-and-Pepper Noise:** Randomly replace 1% of the image pixels with white ("salt" noise) and black ("pepper") pixels.
- **Dirty Rollers:** Add random rollers with thickness between 1 and 3 pixels, a line addition probability of 0.05 per pixel row or column.
- **Random Rotation:** Apply a random rotation of -3° and $+3^\circ$.
- **Binarization:** Utilize the Augraphy⁹ to simulate effects such as dilation, erosion, and letterpress printing.
- **Warping:** Apply geometric transformations and folding effects via Augraphy to mimic paper creases.
- **Shadows:** Apply light gradients and shadow cast from Augraphy to simulate shadows that occur when a document is taken.
- **Blur via Point Spread Function:** Generated 100 PSF kernels and randomly applied one to the document.

We classify above distortions into two categories: (1) weak distortions: These preserve text clarity and include background addition, binarization, minor rotation, and PSF-based blurring. (2) strong distortions: These degrade readability, causing blurriness and font warping. They include salt-and-pepper noise, dirty rollers, warping, and shadow effects. To simulate varying levels of document distortion, we apply the above strategies in three ways:

- Apply a weak distortion per page.
- Apply a strong distortion per page.
- Apply two randomly selected distortions per page.

We generate three document image datasets with varying noise levels and parse structured data using MinerU, GOT, and Qwen2.5-VL, resulting in 9 perturbed datasets. The examples of distorted documents are shown in Fig. S1. In Sec. 4.3, we evaluate RAG performance on these datasets,

⁹<https://github.com/sparkfish/augraphy>

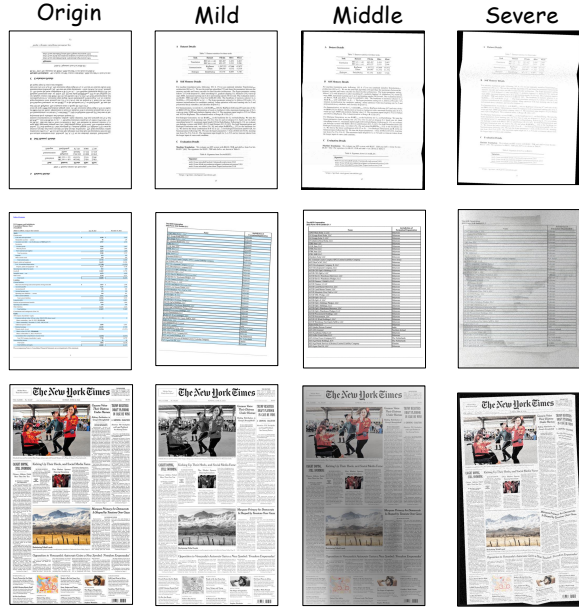


Figure S1. Cases of distorted documents.

reporting the average results for each noise level.

IV. Additional Experimental Results

IV.1. Experimental details

For MinerU, we use version 0.9.2¹⁰ by default. For Marker, version 0.2.17¹¹ is employed. For Nougat, we utilize its 0.1.0-base model (350M). All prompt templates can be found in Sec. I.

For all LLMs and VLMs, we set the temperature to 0 with `do_sample=False` by default for reproducibility.

IV.2. Details in different domains

V. Case Study

Fig. S2 to Fig. S11 show some cases of GOT, MinerU, and Qwen2.5VL-72B on OHRbench. For each case, we indicate the evidence source and answer, giving the OCR result of different models and the responses at the retrieval and generation stages.

¹⁰https://github.com/opendatalab/MinerU/releases/tag/magic_pdf-0.9.2-released

¹¹<https://github.com/VikParuchuri/marker/releases/tag/v0.2.17>

Domain	GT	MinerU	Marker	GOT	Nougat	Qwen2.5-VL	InternVL2.5
Law	81.2	71.0	77.1	62.1	69.0	76.4	69.6
Finance	59.7	36.4	45.0	30.4	25.8	47.9	47.1
Textbook	73.2	43.8	49.6	48.8	37.1	58.3	55.0
Manual	79.1	60.4	68.6	58.9	47.8	71.3	70.2
Newspaper	40.5	31.3	34.0	12.4	10.6	27.7	18.4
Academic	75.1	50.3	55.2	50.2	45.0	61.1	57.1
Administration	82.2	59.4	68.3	57.7	52.7	73.1	73.8
All	70.0	50.1	56.6	45.4	40.8	59.2	55.8

Table S3. Retrieval performance across different domains.

Domain	GT	MinerU	Marker	GOT	Nougat	Qwen2.5-VL	InternVL2.5
Law	56.9	53.4	54.4	43.3	48.8	53.9	50.9
Finance	43.1	30.1	29.5	19.7	17.7	35.9	36.8
Textbook	37.6	25.9	28.2	24.8	16.8	29.1	29.1
Manual	50.2	45.3	46.1	41.3	34.3	48.7	47.7
Newspaper	35.0	33.7	31.6	9.5	8.4	19.6	11.7
Academic	38.3	29.5	27.9	25.3	24.8	33.2	31.3
Administration	46.4	35.7	37.7	32.2	29.2	42.7	42.9
All	43.9	36.1	36.3	27.8	25.5	37.5	35.8

Table S4. Generation performance across different domains.

Domain	GT	MinerU	Marker	GOT	Nougat	Qwen2.5-VL	InternVL2.5
Law	49.6	48.1	48.1	41.1	43.9	47.2	44.9
Finance	27.2	19.4	20.1	15.1	13.1	22.9	22.8
Textbook	30.5	20.9	22.5	21.0	15.7	23.8	23.5
Manual	44.4	38.1	39.8	36.0	30.7	42.3	41.6
News	29.0	25.6	24.7	8.3	5.6	17.4	11.0
Academic	31.9	25.6	24.1	22.8	21.2	27.6	26.4
Administration	41.0	30.9	32.7	29.2	26.6	37.3	37.5
All	36.1	29.5	30.0	24.6	22.2	31.1	29.6

Table S5. Overall performance across different domains.

System:

You are an AI specialized in generating QAs from documents. Your mission is to analyze the document, follow the instructions, and generate RAG-style question-answer pairs based on the document.

RAG-style refers to a question that needs to be answered by retrieving relevant context from an external document based on the question, so the question **MUST** obey the following criteria:

1. Question should represent a plausible inquiry that a person (who has not seen the page) might ask about the information uniquely presented on this page. The questions should not reference this specific page directly (by page number, pointing to a specific paragraph or figure, and never refer to the document using phrases like 'in the document'), nor should they quote the text verbatim. They should use natural language reflecting how someone might inquire about the page's content without direct access.
2. Question must contain all information and context/background necessary to answer without the document. Do not include phrases like "according to the document" in the question.
3. Question must not contain any ambiguous references, such as 'he', 'she', 'it', 'the report', 'the paper', and 'the document'. You **MUST** use their complete names.

User:

Your task is to generate several RAG-style question-answer pairs with different levels of difficulty and evidence sources. {detailed_task_description}.

You **MUST** obey the following criteria:

- The question **MUST** be detailed and be based explicitly on information in the document.
- The question **MUST** include at least one entity.
- The context sentence the question is based on **MUST** include the name of the entity. For example, an unacceptable context is "He won a bronze medal in the 4 × 100 m relay". An acceptable context is "Nils Sandström was a Swedish sprinter who competed at the 1920 Summer Olympics."
- The answer form should be as diverse as possible, including [Yes/No, Numeric, String, List].
- {additional_task_criteria}

If there are no possible questions that meet these criteria, return 'None' as the question. Output the question in JSON format.

{qa_examples}

<document>{document}</document>

Table S6. Q&A Generation Prompt

Structure data task:

In the given documents, the chart elements are all enclosed within `<chart>` `</chart>` tags and illustrated in LaTeX table format. Pay attention to the difference between them and tabular data, as tabular data is not enclosed by `<chart>` `</chart>` tags. **# This paragraph is only used for chart data.**

In order to generate this type of question-answer pairs, first, you need to read the given document, identify the table/formula/chart elements within it, and use them as the evidence context. The evidence context can be a single paragraph for single-hop questions, or several related paragraphs for generating multi-hop questions that require reasoning. After that, you need to generate questions and corresponding answers based on them.

Reading order task:

Your task is to generate RAG-style question-answer pairs from the given two documents. In order to generate this type of question-answer pairs, first, you need to read the given two documents (A, B), identify the text sharing the same entities, and design a question-answer pair based on the contents of both documents A and B. If it is based on the message of document A or document B alone, it cannot be answered.

Understanding task:

You should generate question-answering pairs that require the responder to extract information from documents. The answer should be able to find directly in the documents without any reasoning.

Reasoning task:

You should generate question-answering pairs that require responder to reason before answering, such as calculations, comparisons, finding the maximum and minimum, or integration information from different parts of the documents. The answer should not be able to be found directly in the documents.

Table S7. Detailed description used to generate Q&A pairs for different tasks.

Multi-page Q&A from single-page question:

Your mission is to generate RAG-style combined questions from two questions that have the same entity.

When generating a combined question, there are some criteria you should follow:

- The answer to the combined question should be the same as the answer2.
- It must combine the answer1 to question1 to answer the combined questions. This means that, to answer the combined question, a responder must first deduce the part of the combined question that refers to the answer1, and then proceed to answer the combined question based on that answer.
- You cannot include the answer to question 1 in the combined question.

[{combined_qa_examples}](#)

Based on the above 3 examples, provide a combined question for the following case. If you find it is hard to create such a combined question, output None as the answer. Enclose the combined question within <answer></answer>:

question1: {q1}

answer1: {a1}

question2: {q2}

answer2: {a2}

Multi-page Q&A from different paragraphs:

Your task is to generate RAG-style question-answer pairs from the given two documents and entity names. The entity names appear in both documents, and you need to use them as a bridge to generate the RAG-style question-answer pairs that need to be answered by combining information from both documents.

To generate the question-answer pairs, first, you need to read the given two documents (A, B) and the entity names, find paragraphs related to them, use the paragraphs as evidence context, and design a question-answer pair based on the evidence context from the two documents.

Table S8. Detailed description used to generate multi-page Q&A pairs from both single-page questions and different paragraphs sharing same entities.

System:

You are an expert, you have been provided with a question and documents retrieved based on that question. Your task is to search the content and answer these questions using the retrieved information.

You **MUST** answer the questions briefly with one or two words or very short sentences, devoid of additional elaborations.

Write the answers within <response></response>.

User:

Question: {question}

Retrieved Documents: {retrieved_documents}

Table S9. LLMs prompt for RAG generation

Simple Prompt:

Please do OCR on the image and give all the text content in markdown format. The formulas should be wrapped in \$\$\$. The table and charts should be parsed in LaTeX format. Only output the OCR results without any extra explanations or comments.

Table S10. Simple prompt for VLMs OCR

Detailed Prompt:

You are a powerful OCR assistant tasked with converting PDF images to the Markdown format. You MUST obey the following criteria:

1. Plain text processing:

- Accurately recognize all text content in the PDF image without guessing or inferring.
- Precisely recognize all text in the PDF image without making assumptions in the Markdown format.
- Maintain the original document structure, including headings, paragraphs, lists, etc.

2. Formula Processing:

- Convert all formulas to LaTeX.
- Enclose inline formulas with $$. For example: This is an inline formula $E = mc^2$.$
- Enclose block formulas with $$. For example:
$$\frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$
.$

3. Table Processing:

- Convert all tables to LaTeX format.
- Enclose the tabular data with \begin{table} \end{table} .

3. Chart Processing:

- Convert all Charts to LaTeX format.
- Enclose the chart data in tabular with \begin{table} \end{table} .

4. Figure Handling:

- Ignore figures from the PDF image; do not describe or convert images.

5. Output Format:

- Ensure the Markdown output has a clear structure with appropriate line breaks.
- Maintain the original layout and format as closely as possible.

Please strictly follow these guidelines to ensure accuracy and consistency in the conversion. Your task is to accurately convert the content of the PDF image using these format requirements without adding any extra explanations or comments.

Table S11. Complex prompt for VLMs OCR

System:

You are an AI specialized in document question-answering verification. Your mission is to analyze the given question-answering pairs and follow the instructions. Your response must be true and accurate, and no additional content should be output.

1. Question type check

Dose the question match the task description: {detailed_task_description}

Make sure the question meets the required task context.

2. Evidence relevance Check

Dose the provided evidence context relate to the question provided? Does the answer accurately reflect the information in the evidence context? Ensure the question is formulated based on information explicitly stated. The question should not introduce concepts unrelated to the document's content.

3. Clarity and Precision

Is the question clear and unambiguous? And is the answer concise and precise? Ensure the language is straightforward and easily understandable, and avoid complex phrasing that may confuse the reader. The intention of the question and answer pair must be clear and direct, avoiding verbosity and unnecessary detail. Ensure the answer fully addresses the question without omitting crucial information.

{qas}

Table S12. Q&A Verification Prompt

Evidence Source: Text

2

who had been vaccinated three months before death, from the arm of a healthy child. Three other children vaccinated at the same time, from the same source, took no hurt. On the eighth day after vaccination, a papular and vesicular rash appeared over the trunk, which rapidly assumed a sloughing character. The eruption was at first taken for small-pox, and when death took place, a fortnight later, an inquest was held on the case, for it was then thought to be syphilis. But Mr. HUTCHINSON pointed out that its evolution as well as its character were not those of syphilitic infection, and he considered it to be a true case of vaccinia passing on to a gangrenous condition—a condition he had sometimes observed to take place in variola. The vaccine marks on the arm were natural.—MR. JONATHAN HUTCHINSON, F.R.C.S., &c. *The Lancet*, December 13th, 1879, p. 873.

"HAD TO BE SUSPENDED."

JENNER had already observed phlegmonous erysipelas to follow vaccination. And in the Foundling Hospital at Petersburg, of 5,000 vaccinated infants, 57 took erysipelas in consequence of vaccination and died. Also in Boston, erysipelas has been seen to follow upon vaccination; and on various occasions vaccination has had to be suspended.—DR. C. SPINZIG, St. Louis, U.S. 1880.

"AFRAID OF THE VIRUS."

One pernicious practice, is that of *vaccinating* children. I am as afraid of the virus I am of the scallage. Not only scallage, but other impurities have been conveyed to the blood of healthy persons by means of vaccination. No one, who has given the subject the least consideration, doubts that impurities may be conveyed in this manner.—DR. A. G. SPRINGSTEEN, Cleveland, Ohio. *Am. Med. Medical Tribune*, January, 1880.

"ITS GHASTLY RISKS."

There has fallen an ugly blot. It is too certain that one objection really formidable does exist—that the operation may, in some few instances, import to the subject of it the poison of a hateful and destructive disease (syphilis), peculiar to the human species, and the fruit and nervous of its virus. On this distasteful subject I shall simply appeal to the printed testimony of MR. JONATHAN HUTCHINSON. . . . Such facts as he has demonstrated, con-

Q: What were the consequences of erysipelas following vaccination in the Foundling Hospital at Petersburg?

"HAD TO BE SUSPENDED."

JENNER had already observed phlegmonous erysipelas to follow vaccination. And in the Foundling Hospital at Petersburg, of 5,000 vaccinated infants, 57 took erysipelas in consequence of vaccination and died. Also in Boston, erysipelas has been seen to follow upon vaccination; and on various occasions vaccination has had to be suspended.—DR. C. SPINZIG, St. Louis, U.S. 1880.

Evidence: (Page 3) And in the Foundling Hospital at Petersburg, of 5,000 vaccinated infants, 57 took erysipelas in consequence of vaccination and died.

A: 57 vaccinated infants died

GOT

- ✓ [OCR]: ... JENNER had already observed phlegmonous erysipelas to follow vaccination. And in the Foundling Hospital at Petersburg, of 5,000 vaccinated infants, 57 took erysipelas in consequence of vaccination and died. ...
- ✓ [R]: And in the Foundling Hospital at Petersburg, of 5,000 vaccinated infants, 57 took erysipelas in consequence of vaccination and died.
- ✓ [G]: Consequences: 57 infants died. Vaccination had to be suspended.

MinerU

- ✓ [OCR]: # /n JENNER had already observed phlegmonous erysipelas to follow vaccination. And in the Foundling Hospital at Petersburg, of \$5,000\$ vaccinated infants, 57 took erysipelas in consequence of vaccination and died. ...
- ✓ [R]: of \$5,000\$ vaccinated infants, 57 took erysipelas in consequence of vaccination and died. ...
- ✓ [G]: 57 infants died.

Qwen2.5VL 72B

- ⚠ [OCR]: ...And in the Foundling Hospital at Petersburg, of 5,000 vaccinated infants, 57 took crysipelas in consequence of vaccination and died. Also in Boston,...
- ⚠ [R]: ... 57 took crysipelas in consequence of vaccination and died. Also in Boston...
- ✗ [G]: Consequences: Death, gangrenous condition, ulcerated sores, chronic eczema, debilitation, increased disease, fatal poisoning.

Figure S2. A case using text as the evidence source on a distorted academic document.

Evidence Source: Text

NO.
Date:

2. "A或者B"类议论文模板:
 导入: 第一段: Some people hold the opinion that A is superior to B in many ways. Others, however, argue that B is much better. Personally, I would prefer A because I think A has more advantages.
 正文: 第二段: There are many reasons why I prefer A. The main reason is that... Another reason is that... (赞同A的原因)
 第三段: Of course, B also has advantage to some extent... (列举1~2个B的优势)
 结论: 第四段: But if all these factors are considered, A is much better than B. From what has been discussed above, we may finally draw the conclusion that... (得出结论)
 3. 观点论证类议论文模板:
 导入: 第一段: 提出一种现象或某个观点作为议论的话题
 As a student, I am seriously in favour of the decision (表明自己的观点是赞成还是反对)
 The reasons for this may be listed as follows. (过渡句, 承上启下)
 正文: 第二段: First of all... Secondly... Besides... (列举2~3个赞成或反对的理由)
 结论: 第三段: In conclusion, I believe that... (照应第一段, 构成“总一分一点”结构)
 4. "How to"类议论文模板:
 导入: 第一段: 提出一种现象或某种困难作为议论的话题
 正文: 第二段: Many ways can help to solve this serious problem, but the following may be most effective. First of all... Another way to solve the problem is... Finally... (列举2~3个解决此问题的办法)
 112

Q: In "A或者B"类议论文模版", does the 'A or B' type argumentative essay template conclude by stating that A is better than B?

2. "A或者B"类议论文模板:
 导入: 第一段: Some people hold the opinion that A is superior to B in many ways. Others, however, argue that B is much better. Personally, I would prefer A because I think A has more advantages.
 正文: 第二段: There are many reasons why I prefer A. The main reason is that... Another reason is that... (赞同A的原因)
 第三段: Of course, B also has advantage to some extent... (列举1~2个B的优势)
 结论: 第四段: But if all these factors are considered, A is much better than B. From what has been discussed above, we may finally draw the conclusion that... (得出结论)

Evidence: (Page 0)结论: 第四段: But if all these factors are considered, A is much better than B. From what has been discussed above. We may finally draw the conclusion that: (得出结论)

A: Yes.

GOT

- ! [OCR]: "A或者B"类议论文模板: ... 结论: 第四段: But if all these factors are considered, A is much better than B. From what has been discussed above, we may finally draw the conclusion that. (得出结论)
- ! [R]: ... A is much better than B...
- ✓ [G]: Yes.

MinerU

- ✗ [OCR]: Can't Parse Handwritten Chinese.
- ✗ [R]: Wrong Retrieval Results.
- ✗ [G]: Not answerable

Qwen2.5VL 72B F1-Score:0

- ✓ [OCR]: "A 或者 B" 类议论文模板: ... 结论: 第四段: But if all these factors are considered, A is much better than B. From what has been discussed above, we may finally draw the conclusion that ... (得出结论)
- ✓ [R]: ...considered, A is much better than B....
- ✓ [G]: Yes.

Figure S6. A case using text as the evidence source on a handwritten textbook.

Evidence Source: Text

中考阅读理解备考笔记	
阅读训练的大钟方法	80
1. 讲弯阅读方法	80
2. 训练阅读速度	83
3. 改进阅读方式	84
4. 扩大英语阅读量	85
5. 强化阅读训练	85
6. 加强评估指导	86
做阅读理解的五点注意	86
"先题后文"阅读方法说明	87
阅读题目的答题技巧	87
如何应对阅读理解中的细节题和推理题	88
一. 做细节事实题的方法	88
二. 做推理题的方法	88
1. 事实推断	89
2. 指代推断	89
3. 逻辑推断	89
4. 对作者态度和态度的推断	90
如何推测阅读理解题中的生词词义	90
一. 根据定义或解释推测词义	90
二. 根据语境和逻辑进行判断	91
三. 根据并列或同位关系推测词义	91
四. 根据背景知识和常识判断	91
利用主题句解阅读理解题	92
一. 主题句在段落首尾	92
二. 主题句在段落中间	92

Q: 在中考阅读理解备考中, 阅读训练的六种方法包括哪些?

中考阅读理解备考笔记	80
阅读训练的大钟方法	80
1. 讲弯阅读方法	80
2. 训练阅读速度	83
3. 改进阅读方式	84
4. 扩大英语阅读量	85
5. 强化阅读训练	85
6. 加强评估指导	86

Evidence: (Page 0) 阅读训练的六种方法
80\n1. 讲弯阅读方法 80\n2. 训练阅读速度 83\n3. 改进阅读方式 84\n4. 扩大英语阅读量 85\n5. 强化阅读训练 85\n6. 加强评估指导 86

A: 讲弯阅读方法、训练阅读速度、改进阅读方式、扩大英语阅读量、强化阅读训练、加强评估指导

GOT

! [OCR]: 中青阅读理解各考管记\n\\begin{tabular}{|c|c|}\n\\hline 1. 讲弯阅读方法 & 80\n\\hline 2. 讲练阅读速度 & 83\n\\hline 3. 改进阅读方式 & 84\n\\hline 4. 打大英法阅读量 & 85\n\\hline 5. 强化阅读训练 & 85\n\\hline 6. 加强坪估指导\n\\hline\n\\end{tabular}

✗ [R] Wrong Retrieval Results.

✗ [G]: 阅读训练的六种方法包括: 通读全文、读题查读、简练作答、认真核查、答案形式符合提问方式、应对细节题和推理题的方法。

MinerU

✗ [OCR]: Can't Parse Handwritten Chinese.

✗ [R]: Wrong Retrieval Results.

✗ [G]: Not answerable.

Qwen2.5VL 72B

✓ [OCR]: \$\\begin{array}{ll}\\text{中考阅读理解备考笔记}&80\\text{阅读训练的六种方法}&80\\text{训练阅读方法}&80\\text{训练阅读速度}&83\\text{改进阅读方式}&84\\text{扩大英语阅读量}&85\\text{强化阅读训练}&85\\text{加强评估指导}&86\\end{array}

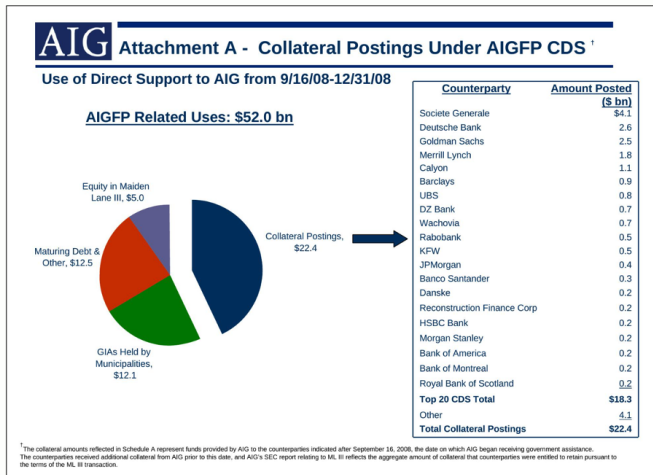
✓ [R]: \$\\begin{array}{ll}\\text{中考阅读理解备考笔记}&80\\text{阅读训练的六种方法}&80\\text{训练阅读方法}&80\\text{训练阅读速度}&83\\text{改进阅读方式}&84\\text{扩大英语阅读量}&85\\text{强化阅读训练}&85\\text{加强评估指导}&86\\end{array}

✓ [G]: \$\\text{训练阅读方法、训练阅读速度、改进阅读方式、扩大英语阅读量、强化阅读训练、加强评估指导。}

Figure S7. A case using text as the evidence source on a handwritten textbook.

Evidence Source: Table

Q: What is the total dollar value of Collateral Postings under AIG's CDS?



Counterparty	Amount Posted (\$ bn)
Societe Generale	\$4.1
Deutsche Bank	2.6
Goldman Sachs	2.5
Merrill Lynch	1.8
Calyon	1.1
Barclays	0.9
UBS	0.8
DZ Bank	0.7
Wachovia	0.7
Rabobank	0.5
KFW	0.5
JPMorgan	0.4
Banco Santander	0.3
Danske	0.2
Reconstruction Finance Corp	0.2
HSBC Bank	0.2
Morgan Stanley	0.2
Bank of America	0.2
Bank of Montreal	0.2
Royal Bank of Scotland	0.2
Top 20 CDS Total	\$18.3
Other	4.1
Total Collateral Postings	\$22.4

Evidence: (Page 2) Collateral Postings \$ 22.4

A: \$22.4 bn.

GOT

- ✓ [OCR]: \begin{tabular}{|c|c|} \hline Counterparty & Amount Posted \\ \hline Societe Generale & \\$4.1 \\ \hline ... Total Collateral Postings & \\$22.4 \\ \hline \end{tabular}
- ✓ [R]: ... \hline Other & 4.1 \\ \n \hline **Total Collateral Postings & (\mathbf{\\$ 22.4})** \\ \hline ...
- ⚠ [G]: Total Collateral Postings: \$22.4

MinerU

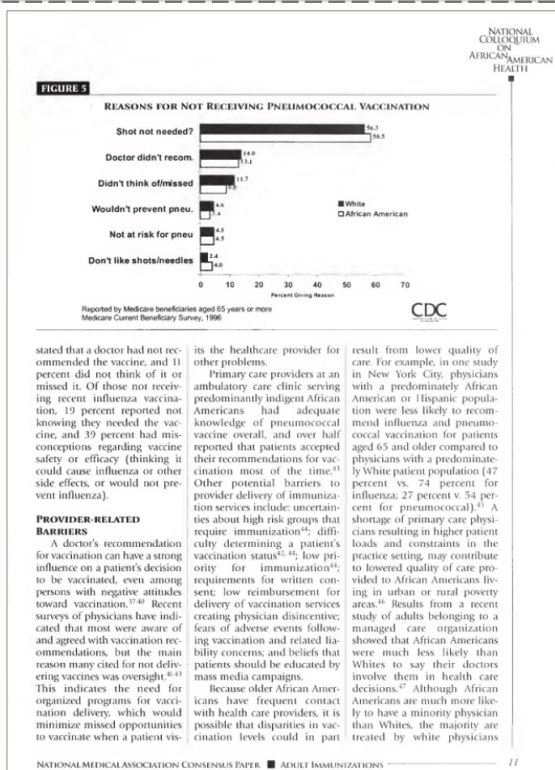
- ✗ [OCR]: Recognize the Table as Chart and Can't Parse Chart.
- ✗ [R]: Content in Other Document.
- ✗ [G]: Not answerable

Qwen2.5VL 72B

- ✓ [OCR]: \begin{tabular}{|c|c|} \hline Counterparty & Amount Posted \\ \hline Societe Generale & \$4.1 \\ \hline ... Total Collateral Postings & \$22.4 \\ \hline \end{tabular}
- ✓ [R]: Other & 4.1 \\ Total Collateral Postings & \$22.4 \\ \n \hline ...
- ✓ [G]: \$22.4 billion

Figure S8. A case using table as the evidence source on a financial report.

Evidence Source: Chart



Q: Among African American respondents, what reason was stated second most for not getting a pneumococcal vaccination?

Evidence: (Page 0)

```
<chart>\n\\begin{tabular}{| c c |}\n\\nReason & White & African\nAmerican \\ \\n\n\\nShot not needed?\n& 56.3 & 58.5 \\ \\n\nDoctor didn't\nrecommend & 14.0 & 13.1 \\ \\n\nDidn't think of/missed\n& 11.7 & 9.0 \\ \\n\nWouldn't prevent pneumonia &\n4.6 & 3.4 \\ \\n\nNot at risk for\npneumonia & 4.5 & 4.5 \\ \\n\nDon't\nlike shots/needles & 2.4 & 4.0 \\ \\n\n\\n\\end{tabular}\\n</chart>"
```

A: Doctor didn't recommend

GOT

✗ [OCR]: **Incorrect Parsed Chart:** \section*{REASONS FOR NOT RECEIVING PNEUMOCOCCAL VACCINATION} \section*{Shot not needed?} Doctor didn't recom. \section*{Didn't think of missed} Wouldn't prevent pneu...

⚠ [R]: \section*{Shot not needed?} \nDoctor didn't recom. \n\section*{Didn't think of missed} \n\n[G]: Didn't think of missed.

MinerU

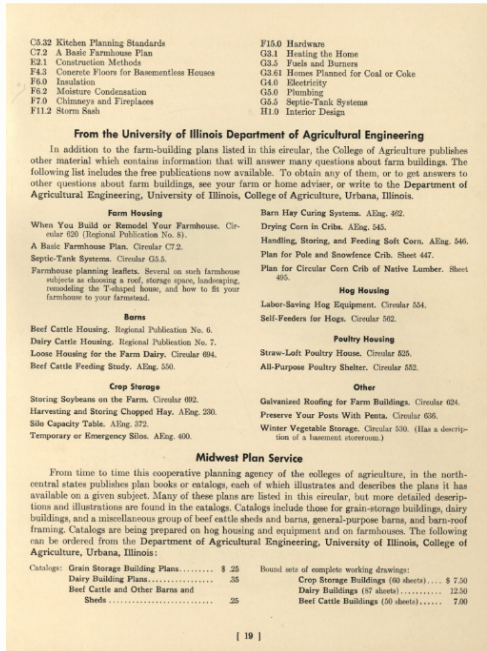
✗ [OCR]: **Can't Parse Chart.**
✗ [R]: REASONS FOR NOT RECEIVING PNEUMOCOCCAL VACCINATION \n ...
✗ [G]: Not knowing they needed the vaccine.

Qwen2.5VL 72B

✓ [OCR]: ... \hline Shot not needed? & \multicolumn{2}{|c|}{56.3} \\ \hline Doctor didn't recom. & \multicolumn{2}{|c|}{14.0} \\ \hline Didn't think of/missed & \multicolumn{2}{|c|}{11.7} \\ \hline ...
✓ [R]: Doctor didn't recom. & \multicolumn{2}{|c|}{14.0} ...
✓ [G]: Doctor didn't recom.

Figure S9. A case using table as the evidence source on a scanned academic paper.

Evidence Source: Text



Q: What is the cost of obtaining the complete set of working drawings for Dairy Buildings from the University of Illinois Department of Agricultural Engineering?

Bound sets of complete working drawings:

Crop Storage Buildings (60 sheets)....	\$ 7.50
Dairy Buildings (87 sheets).....	12.50
Beef Cattle Buildings (50 sheets).....	7.00

Evidence(Page0): Bound sets of complete working drawings: Crop Storage Buildings (60 sheets) ... \\\$ 7.50 Dairy Buildings (87 sheets)..... 12.50 Beef Cattle Buildings (50 sheets) 7.00

A: \$12.50

GOT

- ✓ [OCR]: Crop Storage Buildings (60 sheets) \\\$ 7.50 Dairy Buildings (87 sheets)..... 12.50 Beef Cattle Buildings (50 sheets)..... 7.00.
- ✓ [R]: (page 18) ...Dairy Buildings (87 sheets)..... 12.50
- ✓ [G]: \$12.50

MinerU

- ✓ [OCR]: Bound sets of complete working drawings: Crop Storage Buildings (60 sheets).... \\\$s 7.50 Dairy Buildings (87 sheets)..... 12.50 Beef Cattle Buildings (50 sheets).... 7.00
- ✓ [R]: (page 18) ...Dairy Buildings (87 sheets)..... 12.50...
- ✓ [G]: The cost of obtaining the complete set of working drawings for Dairy Buildings from the University of Illinois Department of Agricultural Engineering is \$12.50.

Qwen2.5VL 72B

- ✓ [OCR]:Crop Storage Buildings (60 sheets).....\$ 7.50 \n Dairy Buildings (87 sheets).....12.50 \n Beef Cattle Buildings (50 sheets).....7.00
- ✓ [R]: (Page 18) Bound sets of complete working drawings:\n Crop Storage Buildings (60 sheets).....\$ 7.50\n Dairy Buildings (87 sheets).....12.50\n Beef Cattle Buildings (50 sheets).....7.00...
- ✗ [G]: Not answerable.

Figure S11. A case using text as the evidence source on a distortion manual.