

# UniDoc-Bench: A Unified Benchmark for Document-Centric Multimodal RAG

Anonymous ACL submission

## Abstract

Multimodal retrieval-augmented Generation (MM-RAG) is a key approach for applying large language models and agents to real-world knowledge bases, yet current evaluations are fragmented—focusing on either text or images in isolation, or simplified multimodal setup, failing to capture document-centric multimodal use cases. In this paper, we introduce UniDoc-Bench, the first large-scale, realistic benchmark for MM-RAG built from 70k real-world PDF pages across 8 domains. Our pipeline extracts and links evidence from text, tables, and figures, then generates 1,600 multimodal QA pairs spanning factual retrieval, comparison, summarization, and logical reasoning queries. To ensure reliability, all of QA pairs are validated and rewritten by multiple human annotators and expert adjudication. UniDoc-Bench supports apples-to-apples comparison across four paradigms — 1) text-only, 2) image-only, 3) *multimodal* text-image fusion and 4) *multimodal* joint retrieval — under a unified protocol with standardized candidate pools, prompts, and evaluation metrics. UniDoc-Bench can also be used to evaluate Visual Question Answering tasks. Our experiments show that multimodal text-image fusion RAG systems outperform both unimodal and jointly multimodal embedding-based retrieval, indicating that neither text nor images alone are sufficient and that current multimodal embeddings remain inadequate. Beyond benchmarking, our analysis reveals when and how visual context complements textual evidence, uncovers systematic failure modes, and offers actionable guidance for developing more robust MM-RAG pipelines.

## 1 Introduction

Retrieval-augmented generation (RAG) has become a widely used approach for applying large language models (LLMs) and agents to real-world knowledge bases (Gao et al., 2023; Fan et al., 2024).

The dominant text-only pipeline applies Optical Character Recognition (OCR) (Li et al., 2022; Xue et al., 2024; Poznanski et al., 2025) to flatten document pages into text, indexes them as chunks, retrieves top-k text passages, and feeds them to a generator. However, many answers depend on information embedded in figures, charts, tables, and complex layouts, where OCR often discards crucial spatial and visual semantics (e.g., map, axes, bar lengths, color encodings) (Ma et al., 2024a; Faysse et al., 2024). These limitations have driven the rapid development of multimodal RAG (MM-RAG), which embeds documents across modalities (text, tables, and images) and retrieves and reasons over them jointly, emerging as a key paradigm for document intelligence.

Current MM-RAG evaluation benchmarks exhibit substantial limitations, as summarized in Table 1. Many are restricted to a single image or a single document page as reference (Mathew et al., 2021, 2022; Zhu et al., 2022; Li et al., 2024; Ma et al., 2024b), cover narrow domains (Mathew et al., 2021, 2022; Zhu et al., 2022; Li et al., 2024), under-represent modalities (Li et al., 2024; Mathew et al., 2022), operate at limited scale (few queries/pages) (Ma et al., 2024b; Wang et al., 2025a) or lack a highly relevant database for RAG evaluation (Ma et al., 2024b). These gaps hinder fair and comprehensive comparison across methods. Moreover, debatable claims have emerged — such as that “image retrieval is all you need” (Faysse et al., 2024; Su et al., 2025) or that multimodal retrieval is inherently superior (Zhang et al., 2024a; Yu et al., 2024a)— without enough fair and unified evaluation. In response, we introduce UniDoc-Bench, a human-verified benchmark spanning 8 domains and covering text, chart, and table content, explicitly designed for cross-modality grounding with examples shown in Figure 1. Crucially, UniDoc-Bench enables apples-to-apples evaluation of text-retrieval, image-retrieval,

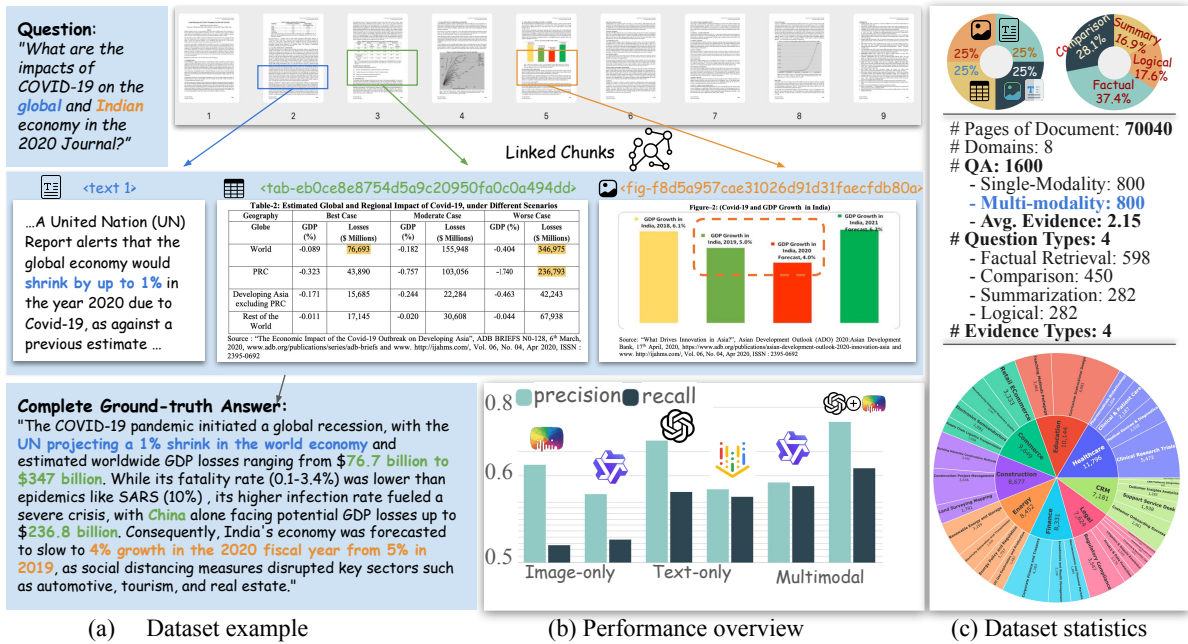


Figure 1: UniDoc-Bench overview.

multimodal text-image-fusion retrieval, and multimodal joint retrieval pipelines using highly relevant large document database and multi-type, cross-modality-grounding queries under a unified protocol. This provides an unbiased view of when multimodal retrieval offers advantages beyond single modalities. In practice, UniDoc-Bench quantifies multimodal gains, guides system design choices, and accelerates the development of effective MM-RAG systems for real-world document intelligence.

We curate a high-quality multimodal RAG evaluation benchmark by designing and applying a classification-based filtering scheme to unlabeled, real-world PDF documents (PDFAs (Montalvo and Wightman, 2024)), yielding 70k highly relevant pages across eight widely used domains—*Finance, Legal, Healthcare, Commerce and Manufacturing, CRM, Energy, Education, and Construction*—containing rich cross-modality content, including text, tables, and images. We construct a knowledge graph that links cross-modality contents across documents via overlapping entities, and leverage these connections to synthesize 1,600 QA pairs spanning four question types: *factual retrieval, comparison, summarization, and logical reasoning*, enabling multi-modality grounding and reflecting realistic retrieval scenarios. To ensure quality, all of the QA pairs are evaluated and rewritten by three independent annotators for faithfulness, completeness, self-containment, human intent, and evidence usability, with disagreements resolved through expert

adjudication. Figure 2 illustrates the full pipeline from PDF segmentation to dataset creation and evaluation.

In this paper, we compare text-only, image-only, multimodal joint, and text-image-fusion retrieval augmented generation pipelines under a unified setup, using identical candidate pools, fixed top- $k$ , consistent prompts, and standardized evaluation criteria. We report retrieval metrics (Recall@10, Precision@10), answer completeness and faithfulness defined at Section 4.2. We observe consistent gains for text-image-fusion RAG systems (completeness = 68.4%) over multimodal joint retrieval systems (64.1%), text-retrieval systems (65.3%), and image-retrieval systems (54.5%). This indicates that retrieving text and images separately using dedicated embeddings, then combining them in the final LLM query, outperforms unified embeddings or single-modality retrieval. Moreover, visual evidence improves answer completeness and enhances faithfulness when paired with textual context, though image-only retrieval cannot fully capture the textual information contained in images. Questions requiring images to answer remain challenging for all systems, suggesting that future RAG improvements should prioritize image-dependent queries. In contrast, performance differences across question types, such as comparison or factual retrieval, are minimal.

We make the following contributions:

Benchmarks	Domain	Evidence	# Queries	# Pages of Doc	RAG Suitable	Unified Evaluation	Multiple Reference	Human Verif
ArxivQA (Li et al., 2024)	single		100k	-	✗	✗	✗	✗
TAT-DQA (Zhu et al., 2022)	single		17k	3k	✗	✗	✗	✓
InfoVQA (Mathew et al., 2022)	single		6k	-	✗	✗	✗	✓
DocVQA (Mathew et al., 2021)	single		11k	-	✗	✗	✗	✓
MMLONG (Ma et al., 2024b)	multiple		1.1k	5k	✗	✗	✓	✓
REALMM (Wasserman et al., 2025)	multiple		5k	8k	✓	✗	✗	✗
ViDoSeek (Wang et al., 2025a)	multiple		1.2k	10k	✓	✗	✗	✗
UniDoc-Bench (ours)	multiple		1.6k	70k	✓	✓	✓	✓

**RAG Suitable:** The dataset provides RAG-style data: queries are self-contained and reflect realistic human questions, with each paired to a grounding corpus (text, images, tables) for retrieval-conditioned answering, supported by a large, highly relevant knowledge base to evaluate retrieval. **Unified Evaluation:** Apples-to-apples comparison across different baseline RAG systems. **Multiple Reference:** Supports multi-hop, multi-modality, multi-source grounding. **Human Verif:** Introduce human experts to review and verify the correctness and quality of all the QA pairs, or to annotate the entire dataset.

Table 1: Comparison of existing document QA datasets with UniDoc-Bench.

- We introduce a new multimodal RAG benchmark built from real-world PDF documents, comprising 70k pages across 8 domains, with 1,600 human-verified QA pairs referencing text, figures, and tables, spanning 4 question types.
- We present a high-quality data synthesizing pipeline for creating MM-RAG evaluation datasets, designed to be compatible with any document database.
- We propose a fair and reproducible evaluation framework by fixing candidate pools across modalities and measuring retrieval effectiveness, answer faithfulness, and completeness end-to-end across different RAG systems.
- We compare text retrieval, image retrieval, text-image fusion, and multimodal joint retrieval pipelines, evaluating which strategy performs best across question types, evidence modalities, and document characteristics. We also show UniDoc-Bench’s use for evaluating Visual Question Answering (VQA) tasks, highlighting its versatility for MM-RAG research.

## 2 Related Works

### 2.1 Multimodal Retrieval-augmented Generation (MM-RAG)

Recent advances in multimodal understanding underscore the importance of MM-RAG for reducing hallucinations. VLM2Vec (Jiang et al., 2024; Meng et al., 2025) shows that instruction-tuning vision-language models improves embeddings for robust text-image alignment. SeBe (Chen et al., 2025) adapts LLaVA-1.5 (Liu et al., 2024) into a retrieval-oriented model that aligns user queries with external knowledge. GME (Zhang et al., 2024b) proposes a unified multimodal embedding capable of text-to-image, image-to-text, and text-to-text retrieval. Uni-Retrieval (Jia et al., 2025) combines VLMs with prompt-tuning to flexibly handle het-

erogeneous queries and modalities. Routing-based methods like UniversalRAG (Yeo et al., 2025) and UniRAG (Sharifmoghaddam et al., 2025) use adaptive query routing to select the best modality and level of granularity.

### 2.2 Visual Document Evaluation

Document understanding with interleaved text and visuals has led to specialized vision-based RAG pipelines (Yu et al., 2024b; Wang et al., 2025b,c) that process document screenshots directly. For example, ColPali (Faysse et al., 2024) uses VLMs to jointly encode textual queries and visual documents via MaxSim (Khattab and Zaharia, 2020), while ViDoRAG (Wang et al., 2025b) employs multi-agent reasoning for iterative cross-modal queries. Optimization-focused methods like VRAG (Wang et al., 2025c) use GRPO (Shao et al., 2024), to adapt VLMs for end-to-end document understanding. However, comparisons with text-only baselines are often unfair, as these baselines ignore non-text modalities. Existing evaluations are also limited: MMLongBench-Doc (Ma et al., 2024c) covers long-context multimodal documents but is poorly suited for retrieval; REAL-MM (Wasserman et al., 2025) and VidoSeek (Wang et al., 2025a) lack cross-page and cross-modal evidence; other benchmarks (Mathew et al., 2021, 2022; Zhu et al., 2022; Li et al., 2024) are narrow in scope, covering single images or pages, limited domains, or small scales (Table 1). To fill these gaps, we introduce UniDoc-Bench, a benchmark designed for practical MM-RAG use cases with multi-page, cross-modal evidence and scalable evaluation.

## 3 Dataset Curation

First, a large-scale, high-quality multi-modal database is needed for evaluating RAG systems, where each document contains content-rich fig-

ures, tables and corresponding textual information. Documents should be domain-specific and exhibit high inter-document similarity to evaluate effective retrieval. The construction of this database is detailed in Section 3.1. Then, we require high-quality query–answer pairs to evaluate the RAG system. Each query is designed to reflect realistic human intent and is written as a self-contained question. The corresponding ground-truth answer must be retrievable solely from the curated database and supported by evidence across multiple modalities. In Section 3.2, we describe our synthetic QA pipeline, and in Section 3.3, we validate dataset quality through human annotation.

### 3.1 Source Document Collection

We use PDFa (Montalvo and Wightman, 2024) as our data source, containing diverse formats (e.g., reports, slides, posters) and covering broad domains, but it lacks tags or labels. Therefore, our first step is data filtering to collect a high-quality database. We design a field scheme (Appendix B.1) that captures key metadata, including domain, subdomain, language, modality (e.g., text, tables, figures), image quality (whether the resolution is clear), and text proportion. This allows us to standardize the data and build a high-quality cross-modality database. As shown in Figure 1 (c), we select 8 domains across industries and define subdomains within each, grouping similar documents. To ensure high inter-document similarity, we retain only documents from 3 – 5 related subdomains containing multiple modalities, yielding on  $\sim 8,000$  pages per domain. The final dataset spans *Legal, Commerce and Manufacturing, Education, Energy, Construction, Finance, Healthcare, and CRM*, with detailed subdomain descriptions in Appendix B.2.

### 3.2 Question and Answer Synthesis Pipeline

As shown in Figure 2, we introduce a data-synthesis pipeline for building multimodal RAG evaluation datasets with high-quality QA pairs, compatible with various document databases.

#### 3.2.1 Evidence Collection

**PDF Parsing.** We parse our curated PDF document database<sup>1</sup> by extracting text chunks, tables, and figures, with the latter two stored separately as image files. Within the parsed text chunk, each image and table is replaced with a unique placeholder tag (e.g., «fig-XXX» or «tab-YYY»), along with its corresponding caption and parsed content

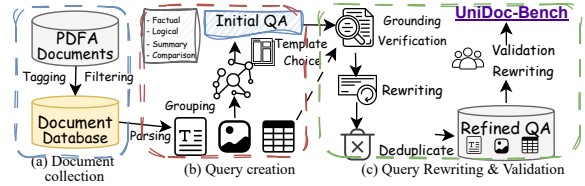


Figure 2: Data Construction pipeline. (a) We filter and tag PDFa documents to curate a high-quality database of 70k pages spanning 8 domains. (b) We parse documents into text, figures, and tables, then synthesize initial QA pairs covering four question types and three modalities using adapted templates. (c) We ground answers in supporting evidence, refine questions for human-intent and self-containment, and verify responses for factuality and completeness, yielding 1,600 QA pairs. To ensure quality, the entire dataset is validated and rewritten by human annotators.

to fully represent interleaved multimodal content. An example is provided in Appendix B.3.

**Chunks Grouping.** To support multimodal evidence QA, we construct a knowledge graph ( $\mathcal{G}_i$ ) (ExplodingGradients, 2024; Peng et al., 2024) over the parsed chunks for domain  $i$ , where nodes ( $N_i = \{n_{i1}, n_{i2}, \dots\}$ ) represent chunks and edges ( $E_i$ ) denote overlapping entities (e.g., “AI Agent Platform”). Chunks across three modalities (text, tables, figures), from within or across documents, are linked to form ground-truth evidence, which are then used for QA synthesis in the next step.

#### 3.2.2 Question and Answer Generation

**Template Choice.** First, we ensure the synthesized questions are **diverse** and span multiple categories, since focusing on a single category or using only the same few-shot example questions can introduce bias and limit the comprehensiveness of RAG evaluation. We designed 4 RAG question types: 1) factual retrieval, 2) comparison, 3) summarization, and 4) logical reasoning. For each type and database domain, we design 10–15 templates (Appendix B.4). We then sample linked chunks ( $n_{ij}, e_{ij}, n_{ik}$ ) and prompt the LLM to select 1–3 templates ( $T_{ij}$ ) that best match the provided chunks and are most likely to produce QA pairs that humans would naturally ask, thereby improving both the diversity and coverage of the questions.

**Evidence Grounding.** To ensure comprehensive evaluation of MM-RAG, we design 4 *answer types* with distinct evidence requirements, each supported by specialized prompts:

- Text-only: The question can be fully answered using natural language text from the documents.
- Image-only: The question requires information

<sup>1</sup><https://unstructured.io/>

exclusively from an image, such as numerical values shown only in a figure.

- **Image-plus-text:** Answering the question requires both text and images, testing the model’s ability to reason across modalities.
- **Table-required:** The question required tabular information to answer, requiring the system to understand table structure and content.

To construct QA pairs, we prompt GPT-4.1 with parsed text chunks and extracted figures/tables (PNG format), guided by prompts  $P_n$  corresponding to the above answer types (see details in Appendix B.5) and templates  $T_{ij}$ . We then employ Gemini-Pro-2.5 — to mitigate single-LLM bias — to verify that the ground-truth answers are correctly grounded in the referenced text, tables, or images, ensuring factual correctness and re-classifying question types when necessary.

**Rewriting.** To ensure that questions are **self-contained** and reflect realistic **human intent**, we refine the initially synthesized QA pairs. In the first stage, many synthesized questions follow a long-context QA style and may include vague references such as “in this report” or “in Figure 8”. To make them suitable for RAG evaluation, we rewrite these questions to ensure they are self-contained and understandable without external context (Appendix B.6). Also, many QA pairs are grounded in images, leading to VQA-style questions (e.g., “How many logos are in Apple Inc.’s 2023 report?”), which do not reflect natural human queries in a RAG context, so we filter and rewrite them to better align with realistic human intent. To ensure comprehensive evaluation, ground-truth answers must be **complete** and **diverse**. In the final step, we revise answers to cover all relevant aspects of their corresponding questions (see Appendix B.7).

### 3.3 Dataset Quality

We evaluate whether our UniDoc-Bench is of sufficient quality to support reliable evaluation of different RAG systems by recruiting 5 human annotators to evaluate the 1,600 question–response pairs against the provided source documents. The annotation process involved assessing each question–response pair across five dimensions (Appendix C):

- **Factuality:** evaluates whether the claims made in the question (Factuality–Question) and the response (Factuality–Response) were factually supported by the source documents.
- **Completeness:** assesses whether the response incorporates all necessary information from the

	Fact.–Q	Fact.–R	Complete.
(%)	98.61	94.20	93.63
	Self-Cont.	Human-like	Grounding
(%)	98.25	96.25	84.38

Table 2: Human evaluation quality on the 1,600 QAs.

- retrieved sources to fully answer the question.
- **Grounding:** assesses whether each source chunk (text, image, or table) used to generate the ground-truth response is required to answer the question, by labeling it as either required or not required, and these labels serve as the ground truth. We then compare the labels produced by our pipeline against the human-annotated ground truth to compute accuracy.
- **Self-Contained:** assesses whether the question was understandable and answerable on its own, without needing external context beyond the provided documents.
- **Human-like Intent:** evaluates whether the question reflected a natural, meaningful query that a human would ask to retrieve information.

As shown in Table 2, the sample shows near-perfect question factuality and self-containment, with strong response factuality and completeness. Human-like intent remains very high (96.25%). Grounding label accuracy is also solid (84.38%). Any questions or responses that do not receive uniformly positive labels are revised by human annotators. These results demonstrate the high quality of UniDoc-Bench for evaluating MM-RAG systems, as well as the robustness of our synthesis pipeline, which can be readily used to generate reliable QA pairs for new databases.

**Dataset Statistics.** UniDoc-Bench consists of 200 QA pairs for each domain, in total 1600 **human-verified and revised** QAs. Within each domain, we have an equal distribution of 50 text-only, image-only, text-plus-image, and table-required questions. In total, the dataset contains 800 single-modality and 800 multi-modality questions. On average, each question requires 2.15 evidence items (text chunks, images, or tables) for a complete answer. More details can be found in Figure 1(b).

## 4 Experiments

To fairly evaluate different RAG systems, we focus on two aspects: retrieval and end-to-end performance. In this section, we first evaluate the retrieval performance of 4 embedding and retrieval models, including text-only, image-only, and two

multimodal approaches (§ 4.1). Next, we evaluate the end-to-end response quality of nine RAG systems that differ in their embeddings, retrieval strategies, and underlying LLMs (Section 4.2). Finally, we demonstrate how our dataset can be used for VQA tasks (Section 4.3). Together, these experiments highlight the usefulness of our dataset and provide practical guidance for selecting RAG components and evaluating VQA systems.

#### 4.1 Retrieval Performance

**Baselines.** We use the curated PDF documents as the knowledge base and the synthesized 1,600 QA pairs to evaluate the following 4 embedding–retrieval models. For all methods, we retrieve the top- $k = 10$  candidates.

- **Text:** PDFs are parsed<sup>1</sup> into text chunks, each embedded with OpenAI’s `text-embedding-3-small`, and retrieved via vector search.
- **Image:** Each PDF page is converted to an image, which is embedded using ColQwen 2.5-v0.2 (Faysse et al., 2024) for image retrieval.
- **MM:** Both text chunks and page-level images are embedded.
  - **MM (GME):** Text and images are jointly embedded using GME-Qwen2-VL-7B-Instruct (Zhang et al., 2024b), enabling multimodal retrieval.
  - **MM (T+I):** A fusion baseline that selects the top-5 candidates from Text and the top-5 from Image retrieval.

**Metrics.** We report Precision@10 and Recall@10 as the retrieval metrics. Since no re-ranker is applied, recall is more informative than nDCG for evaluation. Since we need to evaluate both image and text retrieval, each retrieved text chunk or PDF image-page is mapped back to its original PDF page, and the ground-truth contexts are mapped in the same way. Consequently, a retrieved chunk may span multiple consecutive pages of the source document (e.g., pages 2–3 of document A). A retrieval is considered a true positive if the retrieved text chunk or image-page matches the ground-truth context in both page number and file. This criterion may slightly inflate Recall@10, since partial overlaps (e.g., retrieved pages 1–3 vs. ground-truth pages 3–5, with the answer on page 5) are still treated as correct. However, this approach offers the most practical and fair basis for comparing text and image retrieval. Thus, absolute scores should not be overinterpreted; the key is the relative per-

Domain	Text		Image		MM (GME)		T+I	
	Prec.	Rec.	Prec.	Rec.	Prec.	Rec.	Prec.	Rec.
Com.	.286	.829	.179	.843	.437	.882	<b>.449</b>	<b>.914</b>
Cons.	.246	.762	.159	.792	<b>.429</b>	<b>.864</b>	.422	.853
CRM	.271	.783	.175	.830	<b>.437</b>	.860	.426	<b>.863</b>
Edu	.278	.855	.160	.851	<b>.432</b>	.878	.427	<b>.896</b>
Energy	.239	.706	.148	.718	.366	.723	<b>.374</b>	<b>.746</b>
Fin.	.254	.781	.177	.818	<b>.434</b>	.891	.427	<b>.912</b>
HC	.297	.746	.151	.856	<b>.455</b>	<b>.859</b>	.455	.851
Legal	.312	.861	.178	.861	<b>.462</b>	.883	.458	<b>.903</b>
<b>Avg.</b>	<b>.273</b>	<b>.790</b>	<b>.166</b>	<b>.821</b>	<b>.431</b>	<b>.855</b>	<b>.430</b>	<b>.864</b>

By Question Type								
F.R.	.205	.747	.140	.825	<b>.226</b>	.859	.219	<b>.869</b>
Comp.	.283	.820	.163	.835	<b>.313</b>	.901	.309	<b>.909</b>
Summary	.336	.828	.200	.800	.355	.880	<b>.360</b>	<b>.895</b>
Logical	.365	.820	.201	.813	<b>.386</b>	.870	.382	<b>.882</b>

By Answer Type								
Text-only	.383	.839	.217	.790	<b>.406</b>	.847	.400	<b>.849</b>
Img-only	.081	.724	.092	.878	.097	.909	<b>.097</b>	<b>.920</b>
Text + Img	.336	.847	.190	.824	.350	.888	<b>.351</b>	<b>.908</b>
Table-req.	.291	.752	.163	.791	<b>.326</b>	.851	.316	<b>.861</b>

Table 3: Retrieval performance (Precision@10 / Recall@10) of four RAG systems on 1,600 QA pairs across eight domains (top) and broken down by question and answer types (bottom).

formance differences across methods.

Table 3 summarizes the retrieval performance of four RAG embedding–retrieval models across eight domains, four question types, and four answer types. We observe that **image-based retrieval achieves consistently higher recall but lower precision than text-based retrieval**, as page-image chunks cover more information than individual text chunks. Combining text and image retrieval (T+I) further improves both recall and precision, effectively leveraging the strengths of both modalities. In contrast, multimodal embeddings (GME-Qwen2-VL-7B-Instruct), which encode text and images jointly rather than separately, achieve comparable precision but lower recall, suggesting that current multimodal embeddings still lag behind fusion of unimodal embeddings.

#### 4.2 End-to-End Performance

**Baselines.** We have following six baselines:

- **Image-only RAG:** Each PDF page is converted to a JPEG and retrieved via image embeddings.
  - **Image-only RAG (IMG):** Uses LlamaIndex with colqwen2.5-v0.2 (Faysse et al., 2024) for image retrieval. After retrieval, the question and retrieved images are provided to GPT-4.1 to obtain the final response.
  - **VRAG (Wang et al., 2025d):** a multimodal RAG agent that uses a vision-specific action

480 space — cropping and scaling — to iteratively  
481 extract information from image-formatted  
482 PDF pages in a coarse-to-fine manner. The  
483 embedding model is colqwen2.5-v0.2, and  
484 the final LLM is GPT-4.1.

- 485 • **Text-only RAG (TEXT):** Most multimodal RAG  
486 studies (Wang et al., 2025a; Faysse et al., 2024)  
487 compare only against text-only baselines. For a  
488 fairer comparison, PDF pages are parsed into text  
489 chunks, embedded for retrieval, with associated  
490 images/tables linked back for final responses. In  
491 this baseline, each text chunk is embedded using  
492 text-embedding-3-small and retrieved. The  
493 retrieved text chunks, along with their associated  
494 images, are then fed into GPT-4.1 to generate the  
495 final response.
- 496 • **MM-RAG:** Both parsed text and image-format  
497 page images are embedded and retrieved.
  - 498 – **Multimodal Text-Image-Fusion RAG**  
499 **(T+I):** Retrieves text and images separa-  
500 tely using text-embedding-3-small and  
501 colqwen2.5-v0.2, then combines them for  
502 generation with GPT-4.1. We also evaluate  
503 multiple state-of-the-art LLMs, including  
504 Gemini-pro-2.5, Claude 4.5, and GPT-5.
  - 505 – **Multimodal-joint-Retrieval RAG (MM):**  
506 Uses GME-Qwen2-VL-7B-Instruct (Zhang  
507 et al., 2024b) (MM(G)) or voyage  
508 -multimodal-3 (MM(V)) as a multi-  
509 modal embedding model for both text and  
510 images. Unlike T+I, where text and images  
511 are embedded and retrieved separately, the  
512 text chunks and image-formatted PDF pages  
513 are embedded together, retrieved jointly, and  
514 then fed into GPT-4.1 for the final response.

515 **Metrics.** For **end-to-end** performance, we use  
516 an LLM-based judge to measure faithfulness and  
517 completeness. Specifically, we first ask the LLM to  
518 extract the facts required to answer each question  
519 and then verify whether these facts are grounded  
520 in the ground-truth chunks; this is measured as  
521 faithfulness ( $\uparrow$ ). Next, we ask the LLM to ex-  
522 tract the facts required to answer the question from  
523 the ground-truth answer and then check whether  
524 each fact appears in the system’s response; this is  
525 measured as completeness ( $\uparrow$ ).

526 Table 4 (red background) reports the complete-  
527 ness of responses generated by the six RAG sys-  
528 tems. **Text-only RAG (0.619) substantially out-**  
529 **performs Image-only RAG systems (IMG: 0.527,**  
530 **VRAG: 0.536), highlighting the significant perfor-**  
531 **mance gap between text-based and image-based**

532 retrieval in current RAG architectures. Although  
533 image retrieval achieves higher recall at the re-  
534 trieval stage, this advantage does not translate  
535 into better end-to-end performance, since multi-  
536 modal LLMs (GPT-4.1) are more effective when  
537 processing text and image chunks together rather  
538 than page-level image PDFs alone. In addition,  
539 the low precision of image retrieval makes it  
540 harder for the model to identify the correct in-  
541 formation. The text-image-fusion RAG (T+I)  
542 achieves the best overall performance (0.654)  
543 across eight domains, demonstrating that image-  
544 based PDF representations can effectively comple-  
545 ment text retrieval. Although VRAG leverages  
546 cropping and scaling to enhance image-based re-  
547 trieval (0.536 for VRAG vs. 0.527 for IMG), it still  
548 lags behind the combined T+I approach, under-  
549 scoring the advantage of explicitly integrating both  
550 modalities. Multimodal joint-retrieval RAG sys-  
551 tems (MM (voyage-multimodal-3): 0.637; MM  
552 (GME-Qwen2-VL-7B-Instruct): 0.639) also fall  
553 short of the simple combination of the best text  
554 and image embeddings. This indicates that current  
555 multimodal embedding approaches still have sub-  
556 stantial room for improvement, and that explicitly  
557 **combining separate text and image embeddings**  
558 **remains the most effective strategy** for lever-  
559 aging multimodal documents. More notably, **in**  
560 **some domains—CRM, Education and Legal—**  
561 **multimodal joint RAG performs worse than text-**  
562 **only RAG**, indicating that current multimodal mod-  
563 els still lag behind strong unimodal baselines in  
564 certain domains. These results highlight the impor-  
565 tance of establishing fair baselines and the value of  
566 UniDoc-Bench: multimodal RAG systems should  
567 be benchmarked against strong, balanced baselines  
568 on diverse and high-quality datasets rather than  
569 against overly weak text-only settings.

570 Table 4 (column T+I) compares different state-  
571 of-the-art LLMs used in the Text&Image Retrieval  
572 setting. Claude-4.5-sonnet achieves the best per-  
573 formance across all domains, question types, and  
574 answer types. The table also shows that questions  
575 requiring only text are most effectively handled by  
576 RAG systems with text-embedding. **Questions re-**  
577 **quiring tables are also relatively easy for RAG**  
578 **systems**, as tables can be accurately parsed as text,  
579 which is a straightforward step before embedding  
580 documents for text-based retrieval. In contrast,  
581 questions requiring images remain challenging  
582 across all embedding types — text, image, or mul-  
583 timodal — highlighting that future **RAG improve-**

Domain	Image-only		Text-only	Multimodal							VQA			GT
	IMG	VRAG	TEXT	MM (V)	MM (G)	T+I				GPT-5	GPT-5	GPT-5	Claude	
	GPT-4.1					GPT-4.1	Gemini	Claude	GPT-5					
Com.	.545	.547	.633	.663	.657	<b>.693</b>	.707	<b>.789</b>	.746	.613/.665	<b>.670/.805</b>	.629/.706	<b>.883</b>	
Cons.	.502	.536	.561	.600	.592	<b>.607</b>	.662	<b>.737</b>	.648	.566/.610	<b>.669/.706</b>	.597/.627	<b>.776</b>	
CRM	.524	.523	.643	.635	.640	<b>.647</b>	.689	<b>.771</b>	.696	.612/.679	<b>.756/.774</b>	.614/.666	<b>.848</b>	
Edu	.569	.517	<b>.692</b>	.660	.673	.688	.672	<b>.765</b>	.637	.612/.636	<b>.720/.741</b>	.612/.651	<b>.845</b>	
Energy	.535	.558	.607	<b>.675</b>	.661	.649	.682	<b>.768</b>	.721	.584/.710	<b>.750/.799</b>	.646/.679	<b>.830</b>	
Fin.	.500	.529	.584	<b>.641</b>	.638	.638	.672	<b>.788</b>	.693	.585/.635	<b>.727/.808</b>	.631/.670	<b>.835</b>	
HC	.481	.481	.602	.628	<b>.651</b>	.621	.689	<b>.767</b>	.665	.580/.673	<b>.723/.735</b>	.604/.647	<b>.849</b>	
Legal	.558	.599	.629	.597	.600	<b>.689</b>	.705	<b>.770</b>	.714	.636/.654	<b>.647/.740</b>	.671/.680	<b>.858</b>	
<b>Avg.</b>	<b>.527</b>	<b>.536</b>	<b>.619</b>	<b>.637</b>	<b>.639</b>	<b>.654</b>	<b>.685</b>	<b>.770</b>	<b>.690</b>	<b>.599/.658</b>	<b>.708/.763</b>	<b>.625/.666</b>	<b>.840</b>	

By Question Type													
F.R.	.557	.344	.648	.619	.612	<b>.677</b>	.687	<b>.739</b>	.694	.569/.685	<b>.709/.756</b>	.645/.698	<b>.829</b>
Comp.	.542	.418	.633	.638	<b>.646</b>	.641	.700	<b>.792</b>	.683	.516/.660	<b>.722/.768</b>	.638/.662	<b>.825</b>
Summary	.536	.407	.626	<b>.652</b>	.649	.640	.666	<b>.759</b>	.689	.530/.638	<b>.670/.777</b>	.596/.627	<b>.867</b>
Logical	.548	.513	.637	.664	<b>.681</b>	.630	.681	<b>.813</b>	.706	.514/.602	<b>.719/.774</b>	.607/.639	<b>.864</b>

By Answer Type													
Text-only	.588	.464	.700	<b>.777</b>	.771	.695	.767	<b>.863</b>	.773	.582/.680	<b>.752/.823</b>	.676/.691	<b>.923</b>
Img-only	.486	.336	.616	.465	.462	<b>.619</b>	.588	<b>.644</b>	.629	.510/.613	<b>.577/.651</b>	.546/.611	<b>.756</b>
Text+Img	.502	.453	.600	.584	.580	<b>.617</b>	.611	<b>.719</b>	.617	.441/.587	<b>.677/.729</b>	.578/.609	<b>.813</b>
Table-req.	.610	.392	.633	.723	<b>.742</b>	.683	.773	<b>.853</b>	.741	.609/.752	<b>.812/.851</b>	.700/.752	<b>.870</b>

Table 4: Completeness of systems evaluated on 1,600 QA pairs across 8 domains. Average recall is reported over all domains, with similarity top- $k$  set to 10. Gemini refers to Gemini-2.5-pro. Claude refers to Claude-4.5-sonnet. For VQA, the first value uses the entire document as image input, while the second value uses ground-truth images only. GT is the performance of Claude-4.5-sonnet on the ground-truth text chunks, images and tables.

ments should prioritize image-required questions. We further observe that multimodal joint RAG achieves stronger performance on text-dominant questions, whereas the T+I RAG is more effective for image-dominant queries. We also provide detailed case studies in Appendix D.

### 4.3 Visual Question Answering Performance

UniDoc-Bench can also be used to evaluate Visual Question Answering (VQA) tasks. Table 4 (gray background) reports the performance of state-of-the-art LLMs — Gemini-pro-2.5, Claude-4.5-Sonnet, and GPT-5 — when applied to entire image-format PDFs and to ground-truth pages only. The results show that Claude-4.5-Sonnet consistently achieves the highest completeness scores across all domains and question types in the VQA setting. All models exhibit a performance gap between the two settings, confirming that reasoning over entire documents is more challenging than over isolated ground-truth images. Gemini-pro-2.5 is the most sensitive to this noise. In contrast, Claude-4.5-Sonnet and GPT-5 are more robust to full-document inputs, showing smaller performance drops.

**Additional Findings.** We show the best performance of Claude-4.5-sonnet on the ground-truth chunks in “GT” column of Table 4. Cost and la-

tency are reported in Appendix E.1. Case studies on the impact of content-rich images are presented in Appendix F.1. Analyses of how question type affects difficulty are provided in Appendix D.1, D.2, and F.2. Finally, Appendix F shows that the number of pages and document formats do not significantly affect MM-RAG performance.

## 5 Conclusion

In this paper, we introduced UniDoc-Bench, a large-scale benchmark for document-centric multimodal RAG, built from 70k real-world PDF pages across 8 domains with 1,600 human-verified QA pairs. Our experiments establish a clear performance hierarchy, showing that **text-image fusion RAG performs the best**, consistently outperforming both joint multimodal (MM) RAG and single-modality RAG systems. This key finding demonstrates that fusing separate, strong retrievers for text and images is currently a more effective strategy than relying on a single joint multimodal embedding or a single modality alone. Our analysis further pinpoints image-dependent queries as the primary challenge for all systems. By providing a standardized platform for fair comparison, UniDoc-Bench serves as a crucial resource to guide the development of more robust and faithful document intelligence systems.

## 638 Limitations

639 UniDoc-Bench relies on LLM-synthesized,  
640 template-based queries which, despite human  
641 verification, may lack the linguistic diversity  
642 and conversational dependency (e.g., multi-turn  
643 follow-ups) characteristic of organic user in-  
644 teractions. The evaluation protocol relies on  
645 assumptions, such as treating page-level retrieval  
646 matches as correct—potentially inflating recall  
647 for dense documents—and explicitly excluding  
648 uncaptioned figures under the assumption they are  
649 non-informative. Furthermore, the benchmark is  
650 currently limited to English-centric documents  
651 across eight specific domains and employs LLM-  
652 based judges for end-to-end metrics, suggesting  
653 that findings may not generalize to low-resource  
654 languages or remain robust against inherent judge  
655 model biases.

## 656 References

657 Boqi Chen, Anuj Khare, Gaurav Kumar, Arjun Akula,  
658 and Pradyumna Narayana. 2025. *Seeing beyond: En-  
659 hancing visual question answering with multi-modal  
660 retrieval*. In *Proceedings of the 31st International  
661 Conference on Computational Linguistics: Industry  
662 Track*, pages 410–421, Abu Dhabi, UAE. Association  
663 for Computational Linguistics.

664 ExplodingGradients. 2024. Ragas: Supercharge your  
665 llm application evaluations. [https://github.com/  
666 explodinggradients/ragas](https://github.com/explodinggradients/ragas).

667 Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang,  
668 Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing  
669 Li. 2024. A survey on rag meeting llms: Towards  
670 retrieval-augmented large language models. In *Pro-  
671 ceedings of the 30th ACM SIGKDD conference on  
672 knowledge discovery and data mining*, pages 6491–  
673 6501.

674 Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Om-  
675 rani, Gautier Viaud, Céline Hudelot, and Pierre  
676 Colombo. 2024. Colpali: Efficient document re-  
677 trieval with vision language models. *arXiv preprint  
678 arXiv:2407.01449*.

679 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,  
680 Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen  
681 Wang, and Haofen Wang. 2023. Retrieval-augmented  
682 generation for large language models: A survey.  
683 *arXiv preprint arXiv:2312.10997*, 2(1).

684 Yanhao Jia, Xinyi Wu, Hao Li, Qinglin Zhang, Yuxiao  
685 Hu, Shuai Zhao, and Wenqi Fan. 2025. Uni-retrieval:  
686 A multi-style retrieval framework for stem’s educa-  
687 tion. *arXiv preprint arXiv:2502.05863*.

688 Ziyang Jiang, Rui Meng, Xinyi Yang, Semih Yavuz,  
689 Yingbo Zhou, and Wenhua Chen. 2024. Vlm2vec:

690 Training vision-language models for massive  
691 multimodal embedding tasks. *arXiv preprint  
692 arXiv:2410.05160*.

693 Omar Khattab and Matei Zaharia. 2020. Colbert: Effi-  
694 cient and effective passage search via contextualized  
695 late interaction over bert. In *Proceedings of the 43rd  
696 International ACM SIGIR conference on research  
697 and development in Information Retrieval*, pages 39–  
698 48.

699 Chenxia Li, Weiwei Liu, Ruoyu Guo, Xiaoting Yin,  
700 Kaitao Jiang, Yongkun Du, Yuning Du, Lingfeng  
701 Zhu, Baohua Lai, Xiaoguang Hu, and 1 others.  
702 2022. Pp-ocrv3: More attempts for the improve-  
703 ment of ultra lightweight ocr system. *arXiv preprint  
704 arXiv:2206.03001*.

705 Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong  
706 Feng, Lingpeng Kong, and Qi Liu. 2024. Multimodal  
707 arxiv: A dataset for improving scientific comprehen-  
708 sion of large vision-language models. *arXiv preprint  
709 arXiv:2403.00231*.

710 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae  
711 Lee. 2024. Improved baselines with visual instruc-  
712 tion tuning. In *Proceedings of the IEEE/CVF con-  
713 ference on computer vision and pattern recognition*,  
714 pages 26296–26306.

715 Xueguang Ma, Sheng-Chieh Lin, Minghan Li, Wenhua  
716 Chen, and Jimmy Lin. 2024a. Unifying multimodal  
717 retrieval via document screenshot embedding. *arXiv  
718 preprint arXiv:2406.11251*.

719 Yubo Ma, Yuhang Zang, Liangyu Chen, Meiqi Chen,  
720 Yizhu Jiao, Xinze Li, Xinyuan Lu, Ziyu Liu, Yan Ma,  
721 Xiaoyi Dong, and 1 others. 2024b. Mmlongbench-  
722 doc: Benchmarking long-context document under-  
723 standing with visualizations. *Advances in Neural  
724 Information Processing Systems*, 37:95963–96010.

725 Yubo Ma, Yuhang Zang, Liangyu Chen, Meiqi Chen,  
726 Yizhu Jiao, Xinze Li, Xinyuan Lu, Ziyu Liu, Yan Ma,  
727 Xiaoyi Dong, and 1 others. 2024c. Mmlongbench-  
728 doc: Benchmarking long-context document under-  
729 standing with visualizations. *arXiv preprint  
730 arXiv:2407.01523*.

731 Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthe-  
732 nis Karatzas, Ernest Valveny, and CV Jawahar. 2022.  
733 Infographicvqa. In *Proceedings of the IEEE/CVF  
734 Winter Conference on Applications of Computer Vi-  
735 sion*, pages 1697–1706.

736 Minesh Mathew, Dimosthenis Karatzas, and CV Jawa-  
737 har. 2021. Docvqa: A dataset for vqa on document  
738 images. In *Proceedings of the IEEE/CVF winter con-  
739 ference on applications of computer vision*, pages  
740 2200–2209.

741 Rui Meng, Ziyang Jiang, Ye Liu, Mingyi Su, Xinyi Yang,  
742 Yuepeng Fu, Can Qin, Zeyuan Chen, Ran Xu, Caim-  
743 ing Xiong, and 1 others. 2025. Vlm2vec-v2: Advanc-  
744 ing multimodal embedding for videos, images, and  
745 visual documents. *arXiv preprint arXiv:2507.04590*.

746	Pablo Montalvo and Ross Wightman. 2024.	Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan,	802
747	<a href="#">pixparse/pdfa-eng-wds</a> [dataset]. Hugging	Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu,	803
748	Face Datasets. Accessed August 2025.	Yutong Dai, Michael S Ryoo, and 1 others. 2024.	804
749	Xiangyu Peng, Prafulla Kumar Choubey, Caiming	xgen-mm (blip-3): A family of open large multi-	805
750	Xiong, and Chien-Sheng Wu. 2024. Unanswerability	modal models. <i>arXiv preprint arXiv:2408.08872</i> .	806
751	evaluation for retrieval augmented generation. <i>arXiv</i>	Woongyeong Yeo, Kangsan Kim, Soyeong Jeong, Jin-	807
752	<i>preprint arXiv:2412.12300</i> .	heon Baek, and Sung Ju Hwang. 2025. Universalrag:	808
753	Jake Poznanski, Aman Rangapur, Jon Borchardt, Jason	Retrieval-augmented generation over multiple cor-	809
754	Dunkelberger, Regan Huff, Daniel Lin, Christopher	pora with diverse modalities and granularities. <i>arXiv</i>	810
755	Wilhelm, Kyle Lo, and Luca Soldaini. 2025. olmocr:	<i>preprint arXiv:2504.20734</i> .	811
756	Unlocking trillions of tokens in pdfs with vision lan-	Shi Yu, Chaoyue Tang, Bokai Xu, Junbo Cui, Jun-	812
757	guage models. <i>arXiv preprint arXiv:2502.18443</i> .	hao Ran, Yukun Yan, Zhenghao Liu, Shuo Wang,	813
758	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	Xu Han, Zhiyuan Liu, and 1 others. 2024a. Vis-	814
759	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	rag: Vision-based retrieval-augmented generation	815
760	Zhang, YK Li, Yang Wu, and 1 others. 2024.	on multi-modality documents. <i>arXiv preprint</i>	816
761	Deepseekmath: Pushing the limits of mathematical	<i>arXiv:2410.10594</i> .	817
762	reasoning in open language models. <i>arXiv preprint</i>	Shi Yu, Chaoyue Tang, Bokai Xu, Junbo Cui, Jun-	818
763	<i>arXiv:2402.03300</i> .	hao Ran, Yukun Yan, Zhenghao Liu, Shuo Wang,	819
764	Sahel Sharifymoghaddam, Shivani Upadhyay, Wenhui	Xu Han, Zhiyuan Liu, and 1 others. 2024b. Vis-	820
765	Chen, and Jimmy Lin. 2025. Unirag: Universal re-	rag: Vision-based retrieval-augmented generation	821
766	trieval augmentation for large vision language mod-	on multi-modality documents. <i>arXiv preprint</i>	822
767	els. In <i>Findings of the Association for Computational</i>	<i>arXiv:2410.10594</i> .	823
768	<i>Linguistics: NAACL 2025</i> , pages 2026–2039.	Xin Zhang, Yanzhao Zhang, Wen Xie, Mingxin Li, Ziqi	824
769	Zhaochen Su, Peng Xia, Hangyu Guo, Zhenhua Liu,	Dai, Dingkun Long, Pengjun Xie, Meishan Zhang,	825
770	Yan Ma, Xiaoye Qu, Jiaqi Liu, Yanshu Li, Kaide	Wenjie Li, and Min Zhang. 2024a. Gme: Improving	826
771	Zeng, Zhengyuan Yang, and 1 others. 2025. Think-	universal multimodal retrieval by multimodal llms.	827
772	ing with images for multimodal reasoning: Founda-	<i>arXiv preprint arXiv:2412.16855</i> .	828
773	tions, methods, and future frontiers. <i>arXiv preprint</i>	Xin Zhang, Yanzhao Zhang, Wen Xie, Mingxin Li, Ziqi	829
774	<i>arXiv:2506.23918</i> .	Dai, Dingkun Long, Pengjun Xie, Meishan Zhang,	830
775	Qiuchen Wang, Ruixue Ding, Zehui Chen, Weiqi Wu,	Wenjie Li, and Min Zhang. 2024b. Gme: Improving	831
776	Shihang Wang, Pengjun Xie, and Feng Zhao. 2025a.	universal multimodal retrieval by multimodal llms.	832
777	Vidorag: Visual document retrieval-augmented gen-	<i>arXiv preprint arXiv:2412.16855</i> .	833
778	eration via dynamic iterative reasoning agents. <i>arXiv</i>	Fengbin Zhu, Wenqiang Lei, Fuli Feng, Chao Wang,	834
779	<i>preprint arXiv:2502.18017</i> .	Haozhou Zhang, and Tat-Seng Chua. 2022. Towards	835
780	Qiuchen Wang, Ruixue Ding, Zehui Chen, Weiqi Wu,	complex document understanding by discrete reason-	836
781	Shihang Wang, Pengjun Xie, and Feng Zhao. 2025b.	ing. In <i>Proceedings of the 30th ACM International</i>	837
782	Vidorag: Visual document retrieval-augmented gen-	<i>Conference on Multimedia</i> , pages 4857–4866.	838
783	eration via dynamic iterative reasoning agents. <i>arXiv</i>		
784	<i>preprint arXiv:2502.18017</i> .		
785	Qiuchen Wang, Ruixue Ding, Yu Zeng, Zehui Chen,		
786	Lin Chen, Shihang Wang, Pengjun Xie, Fei Huang,		
787	and Feng Zhao. 2025c. Vrag-rl: Empower vision-		
788	perception-based rag for visually rich information		
789	understanding via iterative reasoning with reinforce-		
790	ment learning. <i>arXiv preprint arXiv:2505.22019</i> .		
791	Qiuchen Wang, Ruixue Ding, Yu Zeng, Zehui Chen,		
792	Lin Chen, Shihang Wang, Pengjun Xie, Fei Huang,		
793	and Feng Zhao. 2025d. <a href="#">Vrag-rl: Empower vision-</a>		
794	<a href="#">perception-based rag for visually rich information</a>		
795	<a href="#">understanding via iterative reasoning with reinforce-</a>		
796	<a href="#">ment learning</a> . <i>Preprint</i> , arXiv:2505.22019.		
797	Navve Wasserman, Roi Pony, Oshri Naparstek, Adi Raz		
798	Goldfarb, Eli Schwartz, Udi Barzelay, and Leonid		
799	Karlinsky. 2025. Real-mm-rag: A real-world		
800	multi-modal retrieval benchmark. <i>arXiv preprint</i>		
801	<i>arXiv:2502.12342</i> .		

## A The Use of Large Language Models (LLMs)

We used LLMs for three purposes: (i) polishing grammar and improving readability, and (ii) assisting in the evaluation of RAG outputs (iii) synthesizing the QA pairs. All research ideas and analyses were conducted by the authors, who take full responsibility for the content.

## B Dataset Creation Details

### B.1 Document fields

We classify each PDF document into the following fields:

- *domain*: one or more from {Healthcare, Finance, Technology and Software, Commerce and Manufacturing, Marketing, Arts and Entertainment, Government, Legal, Education, Scientific Research and Development, Customer Relationship Management (CRM), others}
- *subdomain*: optional finer-grained categories
- *date*: year or estimated year (e.g., 2005)
- *language*: language of the document (e.g., en)
- *modality*: possible values include {text, table, figure, formula, image, drawing}
- *quality*: parsing confidence, values {easy-parse, hard-parse}
- *format*: one or more from {form, report, notice, paper, slide, poster, book, newspaper, article, textbook, note, webpage, document, record}
- *text\_proportion*: percentage of textual content (e.g., 25%)

As described in Section 3.1, we do not include every domain or subdomain in our benchmark. Instead, we filter the source data and retain eight highly representative domains.

### B.2 Domain Definitions

We classify documents into domains and subdomains, each with a brief description for clarity. These labels are used for tagging. As detailed in Section 3.1, we filter the source data and retain eight highly representative domains rather than including all possible ones.

<b>Domain</b>	<b>Subdomain</b>	<b>Description</b>
Healthcare	Clinical & Patient Care	Direct provider-patient interaction: diagnosis, treatment, and care management.
Healthcare	Pharmaceuticals & Biotechnology	Development and regulation of drugs, vaccines, and biotechnological products (no patient records).
Healthcare	Medical Devices & Diagnostics	Design, production, and regulation of medical equipment and diagnostic tools (no patient records).
Healthcare	Clinical Research & Trials	Controlled studies testing treatments, drugs, or therapies.
Healthcare	Public Health & Policy	Population-level promotion, disease prevention, accessibility (not individual records).
Healthcare	Other Healthcare Topics	Healthcare economics, law, and alternative medicine.
Finance	Investments & Wealth Management	Stock portfolios, retirement planning, mutual funds, hedge funds.
Finance	Insurance & Risk Management	Health, life, auto, property insurance; actuarial analysis.
Finance	Corporate Finance & Treasury	Budgeting, fundraising, M&A, investor relations, corporate structure.
Finance	Personal Finance & FinTech	Budgeting apps, personal loans, P2P lending, digital wallets.
Finance	Real Estate Finance	Mortgages, REITs, valuations, market dynamics.
Finance	Macroeconomics & Financial Markets	Markets, currency, fiscal/monetary policy, global economics.
Finance	Other Finance Topics	Microfinance, Islamic banking, niche financial products.
Technology & Software	Software Engineering & DevOps	Coding, testing, deployment, CI/CD, APIs.
Technology & Software	Cybersecurity & Information Security	Risk management, encryption, compliance, network defense.
Technology & Software	Data Science, AI & Analytics	ML, pipelines, visualization, BI tools.
Technology & Software	HCI & UX	Design, prototyping, accessibility, usability studies.
Technology & Software	Emerging Technologies	AR/VR, quantum computing, IoT, blockchain.
Technology & Software	Other Tech Topics	Legacy systems, databases, systems architecture.
Commerce & Manufacturing	Supply Chain & Logistics	Procurement, warehousing, transportation, inventory.

<b>Domain</b>	<b>Subdomain</b>	<b>Description</b>
Commerce & Manufacturing	Industrial Engineering & Production	Process optimization, quality control, Lean/Six Sigma.
Commerce & Manufacturing	Retail & E-Commerce	Marketplaces, POS systems, consumer engagement.
Commerce & Manufacturing	Trade Policy & Global Commerce	Tariffs, export-import regulation, global trade.
Commerce & Manufacturing	Other Commerce Topics	Business operations, sales, distribution.
Marketing	Digital Marketing & Advertising	Social media, SEO/SEM, online campaigns.
Marketing	Consumer Behavior & Market Research	Surveys, focus groups, data-driven insights.
Marketing	Branding & Corporate Identity	Logo, image, brand value, messaging.
Marketing	Marketing Analytics & Metrics	ROI, attribution models, dashboards.
Marketing	Other Marketing Topics	Public relations, sponsorships, offline campaigns.
Arts & Entertainment	Performing Arts	Music, theater, dance, performance reviews.
Arts & Entertainment	Visual Arts & Design	Painting, sculpture, illustration, graphic design.
Arts & Entertainment	Film, TV & Media Studies	Criticism, production, audience reception.
Arts & Entertainment	Literature & Writing	Fiction, non-fiction, literary analysis.
Arts & Entertainment	Games & Interactive Media	Video games, role-playing, esports.
Arts & Entertainment	Other Arts Topics	Fashion, photography, cultural heritage.
Government	Public Administration & Policy	Bureaucracy, policymaking, implementation.
Government	Law Enforcement & Security	Policing, intelligence, defense, military studies.
Government	International Relations & Diplomacy	Foreign policy, treaties, global governance.
Government	Elections & Governance	Voting, political systems, representation.
Government	Other Government Topics	Civil rights, immigration, taxation.
Legal	Corporate & Business Law	Contracts, mergers, compliance.

<b>Domain</b>	<b>Subdomain</b>	<b>Description</b>
Legal	Criminal & Civil Law	Courts, trials, disputes, legal rights.
Legal	Intellectual Property Law	Copyrights, patents, trademarks.
Legal	International & Comparative Law	Cross-border legal systems, treaties.
Legal	Legal Theory & Jurisprudence	Philosophy of law, frameworks.
Legal	Other Legal Topics	Niche legal issues, regulatory law.
Education	K-12 Education	Curriculum, pedagogy, assessments.
Education	Higher Education & Academia	Universities, research, accreditation.
Education	Online & Distance Learning	MOOCs, e-learning, virtual platforms.
Education	Education Policy & Reform	Accessibility, standards, funding.
Education	Other Education Topics	Lifelong learning, teacher training.
Scientific R&D	Natural Sciences	Physics, chemistry, biology, earth science.
Scientific R&D	Engineering & Applied Sciences	Electrical, mechanical, civil, aerospace.
Scientific R&D	Medical & Life Sciences	Biomedical, genetics, ecology.
Scientific R&D	Computer Science & Computational Fields	Algorithms, theory, AI, networks.
Scientific R&D	Other Science Topics	Interdisciplinary, niche fields.
CRM	Customer Support & Helpdesk	Call centers, chatbots, support tickets.
CRM	Sales & Lead Management	CRM tools, customer tracking, pipelines.
CRM	Customer Analytics & Insights	Segmentation, lifetime value, churn analysis.
CRM	Customer Experience (CX) & Engagement	Feedback, personalization, loyalty programs.
CRM	Other CRM Topics	Partnerships, integrations, omni-channel strategies.

### B.3 Parsing Examples

We use unstructured to parse each PDF into three components: text chunks, images of figures, and images of tables. Since many figures (e.g., signatures or logos) are not informative, we only retain figures that include captions. Figure 3 shows an example of the parsing output, where figures are represented by placeholders such as «fig-XXX» and the parsed text from the figures.

### B.4 Dataset Templates

This is the templates for the domain: finance. We create different templates for different domains, which can be found in our code files in the supplementary materials.

891

892

893

894

895

made the switch, seeking extra return within the fixed income sector, realized just how exacting the toll was on their portfolio. Their unconscious debalancing into the popular bond managers of the day had destabilized their portfolios and destroyed a significant amount of wealth.

In both cases, the slice of the portfolio allocated to the respective asset class (equity value and fixed income, respectively) remained the same, but the risk composition of the slice changed. This experience can be likened to swapping out fresh broccoli for nutrient-light fried veggie sticks. The investor and the aspiring healthy eater aren't quite as diversified in their portfolio or as balanced in their nutrition, respectively, as they think they are.

### Debalancing Today

So let's apply unconscious debalancing to today's investment landscape, namely, to asset allocation funds. In the past 10 years, "outcome-oriented" investment products have experienced rapid growth. The mandate of these products allows the manager to decide the "what and when" of investing within a wide range of asset class exposures. Certainly, the aim of greater diversification is the combination of a better return and

well as from the greater diversification likely (although hardly guaranteed) in the wide opportunity set of out-of-mainstream markets.

To examine the possibility of debalancing in the asset allocation funds, we begin by surveying all the funds in Morningstar's World Allocation and Tactical Allocation categories that have at least a three-year track record. We divide the resulting 117 funds, which compose our sample, into two groups: the "popular" funds (defined as those with net inflows in 2014) and the

the ones investors are pouring money into—have dramatically different risk profiles than the funds investors are exiting. The popular fund, on average, has far more exposure to U.S. equities than its unpopular counterpart, as measured by its trailing three-year beta to the S&P 500 Index over the risk-free rate. When comparing the two categories of asset allocation funds based on an average flow-weighted beta, the difference in exposures is quite stark, as Figure 1 shows. For example, the average flow-weighted beta of U.S. equity for the popular fund is 1.37

### Debalancing Today

So let's apply unconscious debalancing to today's investment landscape, namely, to asset allocation funds. In the past 10 years, "outcome-oriented" investment products have experienced rapid growth. The mandate of these products allows the manager to decide the "what and when" of investing within a wide range of asset class exposures.

Certainly, the aim of greater diversification is the combination of a better return and

g 7% at n 0% e c r e P -50% -33% -55% -100% -83% -150%

U.S. Large EAFE U.S. Bonds TIPS Commodities EM Equity

Source: Research Affiliates, based on data from Bloomberg and Morningstar.

www.researchaffiliates.com

times compared to 0.60 times for the unpopular fund, a difference of 128%. In the diversifying asset classes, such as emerging market equities, the opposite is true: the popular fund's average flow-weighted beta is

0.03 times, over 80% lower than the unpopular fund's 0.20 times.

Investors are presumably relying on these global and tactical allocation funds to provide some degree of diversification and risk reduction to their portfolios, but the popular strategies may in fact be unconsciously

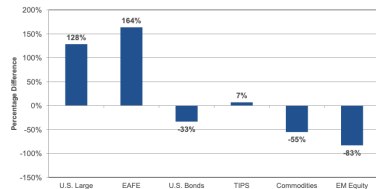
boosting exposure to an expensive asset class with highly unattractive return prospects. By abandoning the unpopular strategies for the popular

ones, investors are unconsciously shifting their risk posture, concentrating their portfolios in the sectors and securities that have

recently outperformed. These securities will inevitably feel the gravitational pull of mean reversion as their valuations

in Figure 2. On average, these strategies don't seem dynamic at all! That's not to say there aren't funds that are tactical and flexible in these categories. There are. But, investors shouldn't expect all of them to be. Our analysis suggests that the popular asset allocation managers, taking

Figure 1. Percentage Difference of Average Flow-Weighted Trailing Three-Year Betas of "Popular" vs. "Unpopular" Funds, as of December 31, 2014



Source: Research Affiliates, based on data from Bloomberg and Morningstar.

Figure 3: Example of PDF parsing with figure placeholders («fig-XXX»).

**Factual Retrieval**

<b>Template</b>	<b>Example</b>
What indicators, policies, or tools are described in the discussion of [Economic Topic/-Financial Strategy]?	What inflation indicators are cited in the ECB's policy blog from June?
Which markets, sectors, or instruments are emphasized in relation to [Trend/Event/Goal]?	Which sectors are favored in the 2025 sustainable investing outlook?
What key positions or exposures are taken by [Investor/Desk/Division] in response to [Condition/Event]?	What position changes did the multi-asset team make in response to rising real yields?
What assumptions, constraints, or parameters are specified in [Scenario/Strategy/Model]?	What assumptions are used in the stress testing scenario for oil price shocks?
When was [Policy/Event/Adjustment] implemented, and what immediate actions followed?	When did the Bank of Japan change its yield curve control stance?
Who oversees or initiates [Financial Decision/Policy/Investment Move] in the described context?	Who approves short-term borrowing requests in the global treasury function?
How is [Strategy/Instrument/Term] defined or operationalized in this context?	How is "duration-neutral tilt" defined in the Q3 fixed income note?
How do you carry out or execute [Action/-Transaction/Plan] in [Financial Context]?	How do you implement a covered call overlay in an income-focused portfolio?
What are the procedural steps or controls listed for [Financial Task/Compliance/Change]?	What steps are required to evaluate bond ladder rollovers in rising rates?

**Comparison**

<b>Template</b>	<b>Example</b>
How do [Strategies/Regions/Instruments] compare in terms of [Risk/Performance/Conditions]?	How do TIPS and gold compare for inflation protection in the current macro setup?
Which asset class, sector, or product is better suited for [Objective/Environment]?	Which is better for income stability in retirement: dividend ETFs or bond ladders?
What are the structural or tactical differences between [Financial Approaches]?	What are the key differences between liability-driven investment and balanced allocation strategies?
How did [Metric/Position/Exposure] change between [Period A] and [Period B]?	How did corporate cash allocation to floating-rate debt shift over 2023?
How do regulatory or monetary responses differ between [Jurisdictions]?	How does Fed liquidity provision compare to ECB emergency facilities post-crisis?

900

**Summarization**

<b>Template</b>	<b>Example</b>
What are the key findings or takeaways from [Brief/Update/Policy/Strategy]?	What are the key points in the tactical asset allocation update from July?
Summarize the main market movements, themes, or risks discussed in [Note/Newsletter/Memo].	Summarize the interest rate risk themes highlighted in the October bond outlook.
What portfolio, liquidity, or policy adjustments are recommended or implemented?	What rebalancing steps were taken in the client model portfolios in Q1?
List the major economic risks or opportunities discussed in [Period/Event/Note].	What macro risks are cited ahead of the U.S. election cycle?
What are the key operational or structural features of [Product/Plan/Tool]?	What are the structural features of the new drawdown facility described in the treasury toolkit?

901

902

**Causal / Reasoning / Why Questions**

<b>Template</b>	<b>Example</b>
Why did [Entity/Desk/Advisor] make [Move/Shift/Decision] in response to [Condition/Event]?	Why did the balanced portfolio reduce international equity in Q2?
How did [Macro Event/Regulatory Shift] influence [Positioning/Allocation/Operations]?	How did the Basel III revisions alter corporate liquidity buffers?
What drove the shift from [Approach A] to [Approach B] in [Context]?	What drove the shift from risk-parity to volatility-targeting in multi-asset allocation?
Why was [Instrument/Policy/Vehicle] introduced or phased out?	Why was the internal netting structure retired in the 2024 treasury overhaul?
What sequence of factors or events led to [Market Reaction/Portfolio Impact/Policy Result]?	What sequence of events led to capital outflows from EM debt in late 2023?

903

## B.5 QA Synthesizing Prompts

### B.5.1 Text-only

#### Prompt P.1: Text-only RAG Question Generation

Prompt: You are an assistant specialized in creating Multimodal RAG tasks. The task is the following: Given some natural language contexts and images inside these contexts, you will generate questions that can be asked by a user to retrieve information from a large documentary corpus.

**Requirements:**

- The 2-hop synthesized question must be a single, self-contained question and must not use "and" to connect multiple questions.
- The answer of the synthesized question will only be found in the contexts.
- The answer of the synthesized question cannot be found in the images.
- The synthesized question must require all the chunks in the contexts to be answered.
- The synthesized question must be specific enough to locate the contexts in a large documentary corpus.
- You must also provide an explanation why the answer can only be found in the provided contexts.

**Question Template:**

- Use the following template to generate the QA:

{{TEMPLATES}}

**Output Format:**

```
{
  "questions": [
    {
      "question": "<synthesized-question>",
      "answer": "<answer-of-the-question>",
      "question_type":
        <choose from "factual_retrieval", "comparison",
          "summarization", "causal_reasoning">,
      "explanation-chunks": "<explanation-chunks>",
      "sentences-chunks-used": {"Chunk1": "sentences-chunk1",
        "Chunk2": "sentences-chunk2", ...}
    }
  ]
}
```

**Input Data:**

- Contexts: "{{contexts}}"
- Images: The image is as follows:

**Notes:**

- If the image can only be used for visualization or illustration, return an empty list for 'sentences-chunks-used'.
- If you cannot use all the chunks in the answer, return an empty list for 'sentences-chunks-used'.

## B.5.2 Image-only

### Prompt P.2: Image-only RAG Question Generation

Prompt: You are an assistant specialized in creating Multimodal RAG tasks. The task is the following: Given some natural language contexts and images inside these contexts, you will generate questions that can be asked by a user to retrieve information from a large documentary corpus.

#### Requirements:

1. The synthesized question must be a single, self-contained question and must not use “and” to connect multiple questions.
2. The answer of the synthesized question will only be found in the image and cannot be found in any sentences in the chunks of the provided contexts.
3. The synthesized question must require chunks/contexts to locate the image and cannot mention the image directly.
4. The synthesized question must be specific enough to locate the contexts in a large documentary corpus.
5. Do not ask “what XYZ in the graph/image/figure”; the question must be general enough to be asked in a large corpus.
6. If you cannot synthesize a question which can only be answered in the image based on the above requirements, do not synthesize anything.
7. Provide an explanation why the answer can only be found in the image and cannot be found in the provided chunks/contexts.
8. Avoid phrasing like “what is shown in the image,” e.g., “what color/logo/name in the image.”
9. Emphasize reasoning, aggregation, temporal comparison, or retrieval from source data. Imagine the question being asked without the image still making partial sense.

#### Question Template:

- Use the following template to generate the QA:

{{TEMPLATES}}

#### Output Format:

```
{
  "questions": [
    {
      "question": "<synthesized-question>",
      "answer": "<answer-of-the-question>",
      "question_type":
        <choose from "factual_retrieval", "comparison",
          "summarization", "causal_reasoning">,
      "image": "<<fig-aaaaa>>",
      "explanation-image": "<explanation-image>",
      "explanation-chunks": "<explanation-chunks>",
      "sentences-chunks-used":
        {"Chunk1": "sentences-chunk1",
         "Chunk2": "sentences-chunk2", ...}
    }
  ]
}
```

#### Input Data:

- Contexts: “{{contexts}}”
- Images: The image is as follows:

#### Notes:

- If the image can only be used for visualization or illustration, return an empty list for ‘sentences-chunks-used’.
- If you cannot use all the chunks in the answer, return an empty list for ‘sentences-chunks-used’.

### B.5.3 Text-plus-Image

#### Prompt P.3: Text-plus-image RAG Question Generation

Prompt: You are an assistant specialized in creating Multimodal RAG tasks. The task is the following: Given some natural language contexts and images inside these contexts, you will generate questions that can be asked by a user to retrieve information from a large documentary corpus.

##### Requirements:

1. The 2-hop synthesized question must require both the provided contexts and images to answer.
2. The concise answer of the synthesized question will directly require information in the image to answer.
3. The concise answer of the synthesized question will also require information in the natural language contexts to answer.
4. The synthesized question must require contexts to locate the image and cannot mention the image directly.
5. The synthesized question must be specific enough to locate the contexts in a large documentary corpus.
6. Provide an explanation indicating which part of the image is used to answer and which sentence in the contexts is used to answer the question.
7. Do not ask “what XYZ in the graph”; the question must be general enough to be asked in a large corpus.
8. If you cannot synthesize a question based on these requirements or directly use the information in the images, do not synthesize anything.
9. If the image can only be used for visualization or illustration, do not synthesize anything. If you cannot use all the chunks in the answer, do not synthesize the question.
10. The synthesized question must be a single, self-contained question and must not use “and” to connect multiple questions.

##### Question Template:

- Use the following template to generate the QA:

```
{{TEMPLATES}}
```

##### Output Format:

```
{
  "questions": [
    {
      "question": "<synthesized-question>",
      "answer": "<answer-of-the-question>",
      "question_type": <choose from "factual_retrieval",
        "comparison", "summarization", "causal_reasoning">,
      "image": "<<fig-aaaaa>>",
      "explanation-image": "<explanation-image>",
      "explanation-chunks": "<explanation-chunks>",
      "sentences-chunks-used":
        {"Chunk1": "sentences-chunk1",
         "Chunk2": "sentences-chunk2", ...}
    }, ...
  ]
}
```

##### Input Data:

- Contexts: “{{contexts}}”
- Images: The image is as follows:

##### Notes:

- If the image can only be used for visualization or illustration, return an empty list for ‘sentences-chunks-used’.
- If you cannot use all the chunks in the answer, return an empty list for ‘sentences-chunks-used’.

## B.5.4 Table-required

### Prompt P.4: Table-required RAG Question Generation

Prompt: You are an assistant specialized in creating Multimodal RAG tasks. The task is the following: Given some natural language contexts containing tables, you will generate questions that can be asked by a user to retrieve information from a large documentary corpus.

#### Requirements:

1. The synthesized question must be a single, self-contained question and must not use “and” to connect multiple questions.
2. The answer of the synthesized question will only be found in the table (within `<table>` and `</table>`) and cannot be found in any sentences outside the `<table>` and `</table>` in the chunks of the provided contexts.
3. The synthesized question must require chunks/contexts to locate the table and cannot mention the ‘table’ directly.
4. The synthesized question must be specific enough to locate the contexts in a large documentary corpus.
5. Do not ask “what XYZ in the table”; the question must be general enough to be asked in a large corpus.
6. If you cannot synthesize a question which can only be answered in the table based on the above requirements, do not synthesize anything.
7. Provide an explanation why the answer can only be found in the table and cannot be found in other parts of the chunks/contexts.
8. Emphasize reasoning, aggregation, temporal comparison, or retrieval from source data. Imagine the question being asked without the table still making partial sense.

#### Question Template:

- Use the following template to generate the QA:

```
{{TEMPLATES}}
```

#### Output Format:

```
{
  "questions": [
    {
      "question": "<synthesized-question>",
      "answer": "<answer-of-the-question>",
      "question_type": <choose from "factual_retrieval", "comparison",
        "summarization", "causal_reasoning">,
      "image": "<<tab-aaaaa>>",
      "explanation-table": "<explanation-table>",
      "explanation-chunks": "<explanation-chunks>",
      "sentences-chunks-used":
        {"Chunk1": "sentences-chunk1",
         "Chunk2": "sentences-chunk2", ...}
    }, ...
  ]
}
```

#### Input Data:

- Contexts: “{{contexts}}”
- Table: The table is included as ‘`<table>`... `</table>`’ in the context.

#### Notes:

- If the table can be used only for visualization or illustration, return an empty list for ‘sentences-chunks-used’.
- If you cannot use all the chunks in the answer, return an empty list for ‘sentences-chunks-used’.

## B.6 Rewriting prompts

### Prompt P.5: Question Rewriting

Prompt: You are tasked with rewriting the following question in two different ways, using only the provided Contexts and without hallucinating any information.

**Date** {{current\_date}}

**Tasks:**

1. **Specific Rewrite:** Add or substitute minimal keywords to tie the question to the Contexts, making retrieval unique while preserving meaning.
2. **Obscured Rewrite:** Paraphrase the specific version to reduce keyword overlap while keeping all needed details intact.

**Requirements:**

- No hallucinated facts.
- Do not remove critical content.
- Avoid source-referencing phrases (“in figure”, “in table”, etc.).
- Rewrites must be standalone, fluent, faithful to Contexts.
- Only add essential keywords (avoid over-specification).

Check if the original answer remains fully correct for both rewrites. If not, set “answer\_wrong” = “True”, else “False”.

**Output Format:**

```
{
  "specific_question":
    "More specific version with essential keywords.",
  "obscured_question":
    "Paraphrased version with reduced keyword overlap.",
  "answer_wrong": "True/False"
}
```

**Example 1:** Original: “What is the revenue growth shown in Figure 3 in 2024’s report?”

```
{
  "specific_question":
    "What is the revenue growth for Company XYZ in 2024?",
  "obscured_question":
    "How did XYZ's financial outcomes change in 2024?",
  "answer_wrong": "False"
}
```

**Example 2:** Original: “What is the median differential rate between hurdle rates and costs of capital for cyclical and non-cyclical firms?”

```
{
  "specific_question":
    "What is the median differential between hurdle
    rates and costs of capital for cyclical vs. non-cyclical firms in
    the S&P 500 according to the Corporate Finance Advisory?",
  "obscured_question":
    "Within the Corporate Finance Advisory, what is the
    median gap between
    required returns and capital costs for S&P 500 firms
    sensitive to the economy vs. stable sectors?",
  "answer_wrong": "False"
}
```

## B.7 Answer Rewriting Prompts

### Prompt P.6: Answer Rewriting

Prompt: You are tasked with rewriting the following answer so that it contains all the facts for answering the question, given the contexts and the image.

**Instruction:**

- Do not hallucinate any additional information. Use only the provided contexts and images.
- The rewritten answer must include the **old correct answer**, if it is correct.
- If the answer is already complete, you may leave it unchanged.
- Make the answer as concise as possible.
- If the **old correct answer** is incomplete, expand it so that the "complete\_answer" fully addresses the question.

**Output Format:**

```
{
  "complete_answer": "Final rewritten answer that is concise,
    faithful to contexts and images, and fully answers the question."
}
```

**Input Data:**

- Question: "{{rewritten\_question\_obscured}}"
- Contexts: "{{contexts}}"
- Old Correct Answer: "{{answer}}"
- Images: The image is as follows:

917 **C Human Annotation**  
918 Annotators were provided with the following in-  
919 structions to evaluate the quality of synthesized  
920 questions and responses against source documents.

### 921 C.1 Task Overview

922 The primary task is to read a synthesized question  
923 and response, then evaluate their quality based on  
924 the provided PDF pages and images. The core  
925 evaluation criterion is factuality.

### 926 C.2 Factuality Evaluation

927 Annotators must determine whether the question  
928 and response are factually supported by the source  
929 material.

#### 930 C.2.1 Procedure

931 Annotators were instructed to follow these steps:

- 932 1. Open the folder corresponding to the given  
933 ID.
- 934 2. Read the text from the PDF pages located in  
935 the chunk\_X subfolder. Annotators were told  
936 to read all text, including tables and image  
937 captions, but to ignore the content of the im-  
938 ages themselves.
- 939 3. Review the images in the img\_X subfolder to  
940 understand which image is being referenced,  
941 then locate that image within the source PDF  
942 to read its context and caption.
- 943 4. Read the provided Question and Response  
944 pair.
- 945 5. Assign a factuality label to both the question  
946 and the response.

#### 947 C.2.2 Label Definitions

948 **Factuality-Question: Factual** All facts and  
949 claims in the question are directly supported  
950 by the source material. There are no  
951 hallucinations or unsupported statements.

952 **Factuality-Question: Not Factual** One or more  
953 facts or claims in the question are not sup-  
954 ported by the source (i.e., contain hallucinated  
955 or fabricated content).

956 **Factuality-Response: Factual** All facts and  
957 claims in the response are directly supported  
958 by the source material. There are no  
959 hallucinations or unsupported statements.

**Factuality-Response: Not Factual** One or more  
960 facts or claims in the response are not sup-  
961 ported by the source (i.e., contain hallucinated  
962 or fabricated content).  
963

**Note:** The original instructions included a rule  
964 stating, "If a question or response is not factual,  
965 it should be labeled as 'Incomplete'." However,  
966 the provided examples use the "Not Factual" label,  
967 which was the standard followed during annotation.  
968

#### 969 C.2.3 Examples

970 The following examples were provided to the an-  
971 notators for guidance.

```
972 {  
973   "id": 0,  
974   "question": "What is the logo of a  
975     major telecommunications company  
976     mentioned in the context related to  
977     personalization strategies?",  
978   "response": "AT&T",  
979 }  
980  
981 # Steps:  
982 # 1. I open folder "0", read all the  
983   chunks and images.  
984 # 2. The question seems factual from one  
985   of the chunk.  
986 # 3. The response seems to NOT be the  
987   correct answer.  
988  
989 # Then, I label Factual-Question as `   
990   Factual`  
991 # Then, I label Factual-Response as `Not  
992   Factual`
```

Listing 1: Example of a factual question with a non-  
factual response.

```
993 {  
994   "id": 4,  
995   "question": "What businesses are  
996     located near the proposed  
997     development area in the Project  
998     Catalyst?",  
999   "response": "AT&T",  
1000 }  
1001  
1002 # Steps:  
1003 # 1. I open folder "4", read all the  
1004   chunks and images.  
1005 # 2. The question seems to be NOT  
1006   factual because I did not see  
1007   Project Catalyst in the pdf or  
1008   images.  
1009 # 3. The response seems to be incorrect  
1010   because the question is not factual.  
1011  
1012 # Then, I label Factual-Question as `Not  
1013   Factual`  
1014 # Then, I label Factual-Response as `Not  
1015   Factual`
```

Listing 2: Example of a non-factual question and  
response.

### 1016 C.3 Completeness Evaluation

1017 This task assesses whether the response provides  
1018 all the necessary information to fully answer the  
1019 question, based on the provided source material.

#### 1020 C.3.1 Procedure

1021 The procedure for evaluating completeness is iden-  
1022 tical to the factuality task: annotators must review  
1023 all provided PDF chunks and images before mak-  
1024 ing a judgment.

#### 1025 C.3.2 Label Definitions

1026 **Complete:** The response includes all the required  
1027 facts and details present in the source material  
1028 needed to comprehensively answer the ques-  
1029 tion.

1030 **Incomplete:** The response omits one or more facts  
1031 or claims that are present in the source and are  
1032 necessary to fully answer the question.

#### 1033 Example 1: Incomplete Response

```
1034 {  
1035   "id": 2,  
1036   "question": "What businesses are  
1037     located near the proposed  
1038     development area in the Project  
1039     Catalyst?",  
1040   "response": "AutoZone Auto Parts,  
1041     Pizza Hut, Sonic Drive In, Joe's  
1042     Pizza Italian",  
1043 }  
1044  
1045 # Steps:  
1046 # 1. I open folder "2", read all the  
1047   chunks and images.  
1048 # 2. The response seems to miss: "Mr Jim  
1049   's Pizza, Justin Spirits, Allsup's  
1050   Convenience Store."  
1051  
1052 # Then, I label Completeness as `  
1053   Incomplete`
```

Listing 3: Example of a response that is missing information available in the source document.

#### 1054 Example 2: Complete Response

```
1055 {  
1056   "id": 0,  
1057   "question": "What is the logo of a  
1058     major telecommunications company  
1059     mentioned in the context related to  
1060     personalization strategies?",  
1061   "response": "AT&T",  
1062 }  
1063  
1064 # Steps:  
1065 # 1. I open folder "0", read all the  
1066   chunks and images.  
1067 # 2. The response seems to be complete.  
1068   AT&T is the only answer.
```

```
# Then, I label Completeness as `  
Complete`
```

Listing 4: Example of a response that contains all necessary information.

### 1072 C.4 Grounding Verification

1073 For each question, annotators were required to  
1074 verify which specific source materials (PDF text  
1075 chunks or images) were necessary to formulate the  
1076 answer.

#### 1077 C.4.1 Procedure and Label Definitions

1078 **Grounding Verification-chunk-X:** After reading  
1079 the question, the annotator must determine  
1080 if the text content of chunk\_X.pdf contains  
1081 any information used in, or required for, the  
1082 answer.

- 1083 • **Required:** The chunk's text contains in-  
1084 formation needed to answer the question.
- 1085 • **Not Required:** The chunk's text does  
1086 not contain any relevant information.

1087 **Grounding Verification-img-X:** The annotator  
1088 must determine if img\_X (including its  
1089 caption and context within the PDF) contains  
1090 any information used in, or required for, the  
1091 answer.

- 1092 • **Required:** The image or its caption con-  
1093 tains information needed to answer the  
1094 question.
- 1095 • **Not Required:** The image and its cap-  
1096 tion do not contain any relevant informa-  
1097 tion.

#### 1098 Example: Grounding Verification

```
1099 {  
1100   "id": 0,  
1101   "question": "What businesses are  
1102     located near the proposed  
1103     development area in the Project  
1104     Catalyst?",  
1105   "response": "AutoZone Auto Parts,  
1106     Pizza Hut, Sonic Drive In, Joe's  
1107     Pizza Italian",  
1108 }  
1109  
1110 # Steps for chunk-0:  
1111 # 1. I open folder "0" and then the sub-  
1112   folder chunk_0.  
1113 # 2. I read the text within pages.pdf.  
1114 # 3. I find part of the answer to the  
1115   question in the text.  
1116 # 4. I label `Grounding Verification-  
1117   chunk-0` as `Required`.  
1118
```

```

1119 # Steps for chunk-1:
1120 # 1. I check for a sub-folder named
1121     chunk_1 in folder "0".
1122 # 2. No chunk_1 sub-folder exists, so I
1123     skip this label.
1124
1125 # Steps for img-0:
1126 # 1. I open folder "0" and then the sub-
1127     folder img_0.
1128 # 2. I view img_0.jpg and locate it in
1129     the original PDF to check its
1130     context.
1131 # 3. I find part of the answer to the
1132     question in the image.
1133 # 4. I label `Grounding Verification-img
1134     -0` as `Required`.
1135
1136 # Steps for img-1:
1137 # 1. I open folder "0" and then the sub-
1138     folder img_1.
1139 # 2. I view img_1.jpg and its context.
1140 # 3. I do NOT find any part of the
1141     answer in this image.
1142 # 4. I label `Grounding Verification-img
1143     -1` as `Not Required`.

```

Listing 5: Example demonstrating how to label individual source chunks and images as required or not required.

## C.5 Self-Contained Evaluation

This task assesses whether a question is understandable and complete on its own, without needing external context or references to specific, unnamed documents.

### C.5.1 Procedure

Annotators were instructed to read only the question and determine if it could be understood and answered without ambiguity, assuming one had access to a large database of documents.

### C.5.2 Label Definitions

**True:** The question is self-contained. It is clearly phrased, makes sense on its own, and provides enough specific detail (e.g., names, topics, concepts) to be answerable. It does not rely on vague document references. For example, "What are the key benefits of solar energy mentioned in the 2022 Department of Energy report?" is self-contained.

**False:** The question depends on external or implicit context to be meaningful. It may contain vague deictic references (e.g., "in the image above," "according to this chart," "what does this mean?") without clarifying what the reference points to. For example, "What is the logo in the image?" is not self-contained as it

requires seeing a specific, un-referenced image.

### Example 1: Not Self-Contained

```

{
  "id": 1,
  "question": "What is the logo in the
  image?",
  "response": "AT&T",
}
# Steps:
# 1. I read the question.
# 2. I find it is NOT clear; "what image
  ?" is an unanswered prerequisite.
# 3. I label `Self-Contained` as `False`.

```

Listing 6: Example of a question that is not self-contained due to a vague reference ("the image").

### Example 2: Self-Contained

```

{
  "id": 0,
  "question": "What is the logo of a
  major telecommunications company
  mentioned in the context related to
  personalization strategies?",
  "response": "AT&T",
}
# Steps:
# 1. I read the question.
# 2. I find it is clear. I can use the
  information within the question to
  search for a relevant document.
# 3. I label `Self-Contained` as `True`.

```

Listing 7: Example of a question that is self-contained because it provides sufficient context ("personalization strategies," "telecommunications company").

## C.6 Human-like Intent Evaluation

This task assesses whether a question reflects a natural and meaningful information-seeking intent, typical of a human user interacting with a document or database.

### C.6.1 Procedure

Annotators were instructed to read the question and judge its authenticity as a genuine human query. The focus was on the nature of the question's intent rather than its grammatical perfection.

### C.6.2 Label Definitions

**True:** The question represents a reasonable and natural query a human would make. It seeks meaningful information such as facts, summaries, comparisons, or explanations, and is

1217 phrased in a way that reflects a real informa-  
1218 tion need. For example: "What were the com-  
1219 pany's main revenue streams in the last fiscal  
1220 year?"

1221 **False:** The question is unnatural, trivial, or does  
1222 not reflect a plausible human intent. This in-  
1223 cludes questions that are overly literal (e.g.,  
1224 counting word occurrences), focus on format-  
1225 ting (e.g., font sizes), are phrased robotically,  
1226 or seek bizarrely specific details that a human  
1227 would be unlikely to ask.

### 1228 **Example 1: Not Human-like**

```
1229 {  
1230   "id": 1,  
1231   "question": "How many logos in the  
1232     Figure one of the major  
1233     telecommunications company?",  
1234   "response": "13",  
1235 }  
1236  
1237 # Steps:  
1238 # 1. I read the question.  
1239 # 2. I do not think a person using an  
1240     information retrieval system would  
1241     ask this style of question.  
1242 # 3. I label `Human-like` as `False`.
```

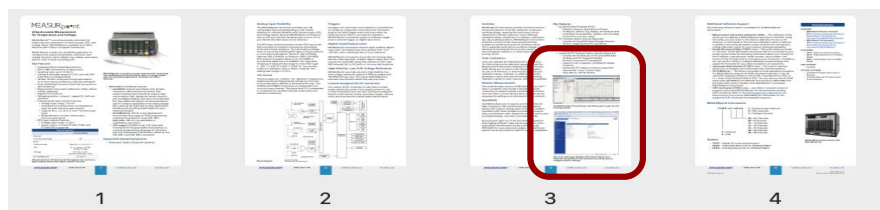
Listing 8: Example of a question that is not human-like due to its trivial, count-based nature.

### 1243 **Example 2: Human-like**

```
1244 {  
1245   "id": 3,  
1246   "question": "What were the top two  
1247     revenues for the EMS division in  
1248     2012?",  
1249   "response": "In 2012, the revenues  
1250     were approximately HK$493,208,000  
1251     and HK$391,677,000.",  
1252 }  
1253  
1254 # Steps:  
1255 # 1. I read the question.  
1256 # 2. I find it is clear and reflects a  
1257     specific, meaningful financial  
1258     inquiry.  
1259 # 3. I label `Human-like` as `True`.
```

Listing 9: Example of a question that reflects a clear, natural, and meaningful information need.

1260	<b>D Examples</b>
1261	<b>D.1 Examples for text-retrieval better than</b>
1262	<b>image-retrieval</b>



**Question:** Which firmware iteration is used by the MEASURpoint device model DT887x-48 from Data Translation Inc.?

**Answer:** The firmware iteration used by the MEASURpoint device model DT887x-48 from Data Translation Inc. is version 1.7.6.3, 6.7.4.86.

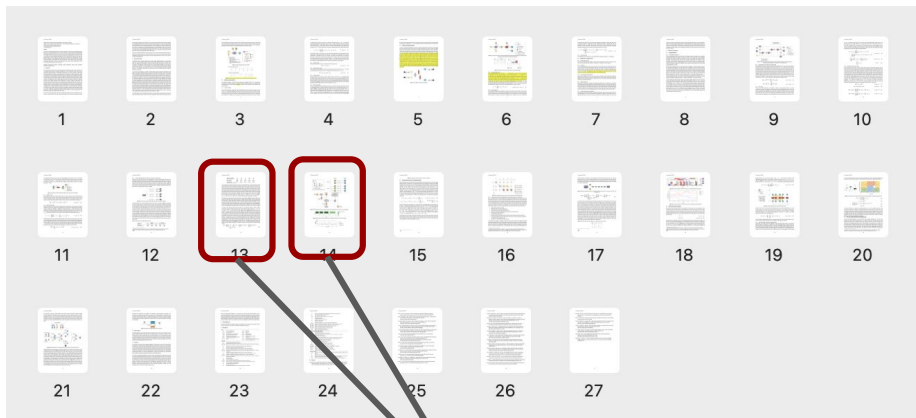
**Text:** The firmware iteration used by the MEASURpoint device model DT887x-48 from Data Translation Inc. is 1.7.6.3, 6.7.4.86. ✓

**Img:** The firmware iteration used by the MEASURpoint device model DT887x-48 from Data Translation Inc. is not explicitly stated in the provided images. The documents mention software and calibration utilities but do not specify a firmware version for the DT887x-48. ✗

**GME:** The firmware iteration used by the MEASURpoint device model DT887x-48 from Data Translation Inc. is displayed by the Eureka Discovery Utility, which provides the version of the firmware running on your instrument. ✗

**T+I:** The MEASURpoint device model DT887x-48 from Data Translation Inc. uses firmware version 1.7.6.3, 6.7.4.86. ✓

Figure 4: Image-retrieval system fails to extract factual facts and details.



**Question:** How many MIbs are produced daily on average for Grade 1 over the span of weeks 0 to 55?

**Answer:** The average daily production level for Grade 1 varies between 179 and 243 MIbs over the span of weeks 0 to 55.

**Text:** The average daily production for Grade 1 over weeks 0 to 55 ranges from **179 MIbs to 243 MIbs**. ✓

**Img:** The average daily production for Grade 1 over the span of weeks 0 to 55 is approximately **2,000 MIbs** (as indicated by the blue "Avg Daily Production (MIbs)" line in Figure 10). ✗

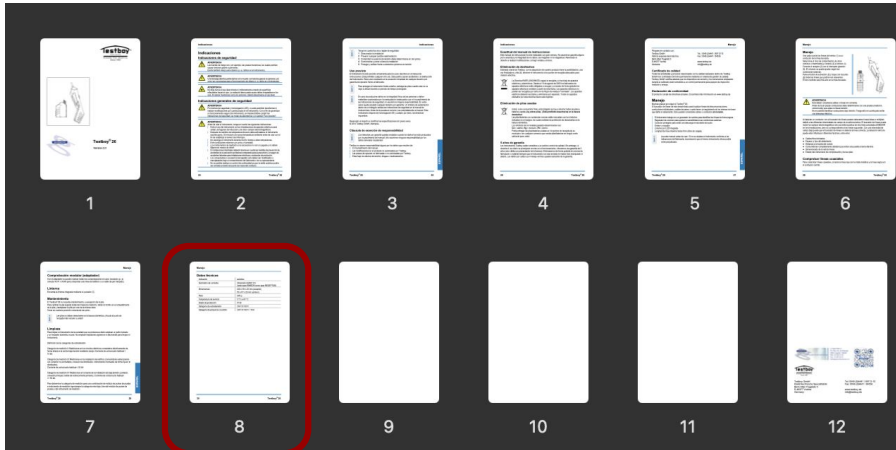
**GME:** Based on Figure 10, the average daily production for Grade 1 over weeks 0 to 55 is approximately **140 MIbs**. ✗

**T+I:** The average daily production for Grade 1 over weeks 0 to 55 ranges from **179 to 243 MIbs**. ✓

Figure 5: Image-retrieval system fails to extract factual facts and details in the image.

1263  
1264

**D.2 Examples for image-retrieval better than  
image-retrieval**



**Question:** Can you provide the service temperature interval for the Testboy® 26?

**Answer:** The service temperature interval for the Testboy® 26 is 5 °C to 40 °C.

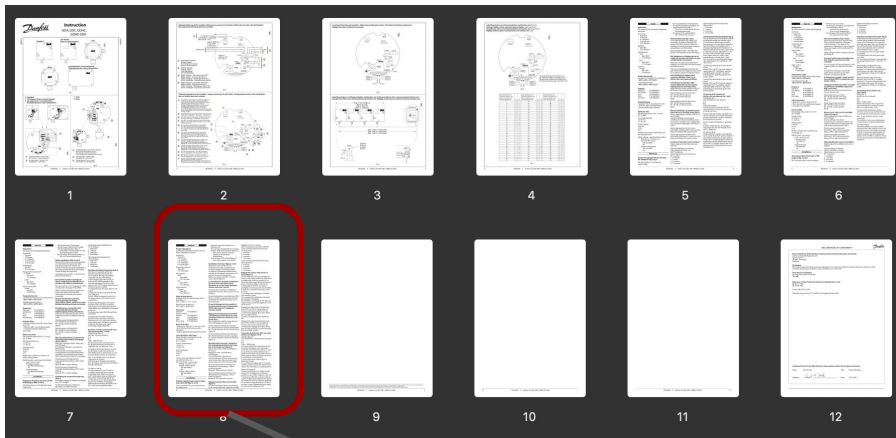
**Text:** The provided contexts do not contain any information about the Testboy® 26 or its service temperature interval. Therefore, I cannot provide the service temperature interval for the Testboy® 26 based on the listed data. ❌

**Img:** The service temperature interval for the Testboy® 26 is 5 °C to 40 °C. ✅

**GME:** The service temperature interval for the Testboy® 26 is 5 °C to 40 °C. ✅

**T+I:** The service temperature interval for the Testboy® 26 is 5 °C to 40 °C. ✅

Figure 6: Text-retrieval system fails to extract factual facts and details in the table.



**Question:** How are the IP protection levels classified for different detector models, such as standard and LCD, used in tough conditions?

**Answer:** Standard and Écran à cristaux liquides detectors have an IP protection level of IP 30 (NEMA 1). Models with IP 65 and EExd have IP 65 (NEMA 4), and the low-temperature model has IP 40 (NEMA 2).

**Text:** The IP protection levels for different detector models, such as standard and LCD, are classified as follows: - Standard and LCD models: No specific IP rating mentioned, but grounding is required for safety. - IP65 models: Protected against dust and low-pressure water jets from any direction (IP65 rating). ❌

**Img:** The IP protection levels for different detector models are classified as follows: - \*\*Standard and LCD display models:\*\* IP 30 (NEMA 1) - \*\*IP 65 model:\*\* IP 65 (NEMA 4) - \*\*EExd model:\*\* IP 65 (NEMA 4) - \*\*Low temperature model:\*\* IP 40 (NEMA 2) ✅

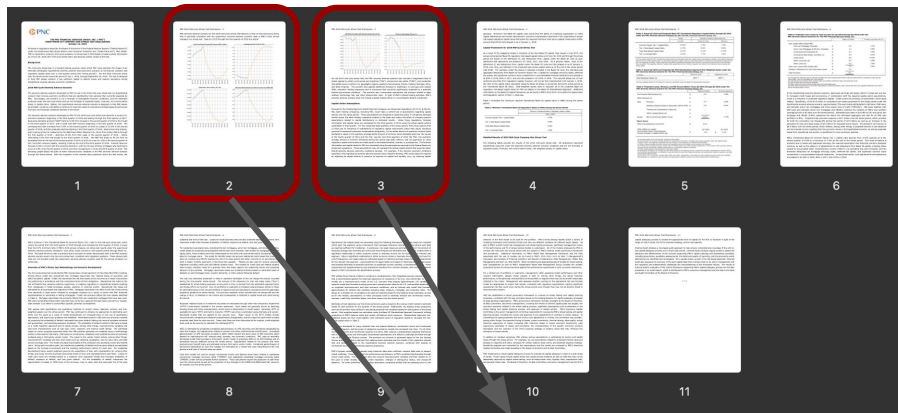
**GME:** The IP protection levels for different detector models are as follows: - Standard and LCD display models: IP 30 (NEMA 1) - IP 65 and EExd models: IP 65 (NEMA 4) - Low temperature model: IP 40 (NEMA 2) For tough conditions, IP 65 and EExd models are used, providing high protection against dust and water jets. ✅

**T+I:** The IP protection levels for different detector models are classified as follows: - Standard and LCD display: IP30 (NEMA 1) - IP 65 model: IP65 (NEMA 4) - EExd model: IP65 (NEMA 4) - Low temperature model: IP40 (NEMA 2) ✅

Figure 7: Text-retrieval system fails to extract factual facts and details in the table.

1265  
1266

**D.3 Examples for multimodal-retrieval better than single-modality-retrieval**



**Question:** According to the PNC's 2016 stress evaluation, how did nonfarm employment levels change between mid-2015 and the third quarter of 2018 in their severe case?

**Answer:** The total nonfarm employment decreased from 144,000 in 2015 Q2 to around 136,000 in 2017 Q1 and then started to recover, reaching approximately 137,000 by 2018 Q3.

**Text:** According to PNC's 2016 stress evaluation, nonfarm employment levels declined from about 144,000 in mid-2015 to around 136,000 by the third quarter of 2018 in their severe case. ❌

**Img:** According to PNC's 2016 stress evaluation, nonfarm employment levels declined sharply from about 142 million in mid-2015 to a low of around 136 million by the third quarter of 2018 in their severe case. ❌

**GME:** According to PNC's 2016 stress evaluation, nonfarm employment levels declined from about 142,000 thousand in mid-2015 to approximately 137,000 thousand by the third quarter of 2018 in their severely adverse scenario. ✅

**T+I:** According to PNC's 2016 stress evaluation, nonfarm employment levels declined from about 144,000 (thousands) in mid-2015 to a low near 136,000 (thousands) before partially recovering to around 137,000 (thousands) by the third quarter of 2018 in their severe case. ✅

Figure 8: MM RAG system handles multi-modality-evidence questions better.

## E Additional Experiments

### E.1 Cost Comparison

**MM-RAG systems can achieve both better end-to-end performance and lower cost than text-only RAG.** Text-only RAG is the most expensive due to high token consumption, while image-only RAG has the lowest cost and latency. Multimodal RAG offers lower cost than text-only RAG with comparable latency.

We report the average inference cost and latency of different RAG systems in Table 6. Image-only systems (IMG) are the most efficient, whereas multimodal systems (MM) are the slowest, reflecting the trade-off between model complexity and capability. The T+I fusion RAG incurs additional latency because it retrieves text chunks before images. Overall, these results show that modern MM-RAG systems can provide improved performance at lower cost than text-only RAG.

## F Additional Analysis

### F.1 Content-rich images increase difficulty

We analyze images from the easiest domains (*commerce manufacturing* and *legal*) and the most challenging domains (*finance* and *construction*). Using gemini-2.5-pro, we classify images as *content-rich* (containing information not present in the text) or *illustrative*. Content-rich images are substantially more common in finance (62.8%) and construction (69.3%) than in commerce manufacturing (40.0%) and legal (49.5%). This suggests that domains with a higher proportion of content-rich images pose greater challenges for RAG, as they require effective multimodal understanding beyond text, consistent with the results in Table 4.

### F.2 Question type affects difficulty

As shown in Section 4.2, the type of context required to answer a question is the most significant factor influencing RAG performance. Different categories of questions contribute unevenly to the advantage of either text- or image-retrieval RAG systems. By carefully analyzing questions that can only be answered correctly by one of the two systems, we summarize the key distinguishing features:

Text-Retrieval Advantages:

- *Entity Recognition* (e.g., brands, organizations; 53.9% of text advantage): Strong at identifying specific people, companies, or organizations.

- *Comparative Analysis* (37.6%): Ranking, evaluating differences, or determining which option is preferable.
  - *Contextual Numerical Reasoning* (34.8%): Numbers requiring understanding of surrounding context.
  - *Quantity Estimation* (29.1%): Questions asking about amounts, counts, or measurements.
  - *Domain-Specific Terminology* (16.3%): Technical, scientific, or specialized terms and standards.
- Image-Retrieval Advantages:
- *Visual Chart Data Interpretation* (64.2% of image wins): Charts and tables make numerical information more accessible. *Example*: How much of the auto ABS senior tranches in Europe were rated AAA in early 2018?
  - *Temporal / Chronological Data* (40.0%): Timeline visualizations clarify temporal relationships. *Example*: When did U.S. petroleum imports drop under \$20 billion?
  - *Technical / Measurement Information* (19.2%): Diagrams often contain measurements or specifications not in text. *Example*: What is the service temperature interval for Testboy<sup>®</sup> 26 based on the listed data?
  - *Spatial / Geographic Reasoning* (13.3%): Maps and layouts convey location context and spatial relationships. *Example*: What is the impact of delivery time on scheduling at 22 Bishopsgate?

### F.3 Document formats do not affect performance.

As discussed in Section 3.1, documents span formats such as newspapers, textbooks, webpages, forms, reports, papers, slides, and posters. In the best-performing domain, *commerce manufacturing*, the distribution is diverse, with reports (45.2%), textbooks (23.6%), papers (18.7%), and webpages (10.5%). In contrast, the worst-performing domain, *finance*, is dominated by reports (80.8%), with only small shares of papers (12.2%), textbooks (2.9%), and webpages (2.3%). Yet this trend is not consistent: the second-worst domain, *construction*, is also diverse, with reports (53.9%), papers (30.4%), and textbooks (11.3%). Therefore, format distribution alone cannot explain performance differences.

**Document layouts do not affect performance.**

In the best-performing domain, *commerce manufacturing*, documents are composed of text (73.9%), tables (4.0%), and figures (22.1%), while the worst-performing domain, *finance*, shows a nearly identical distribution (72.9% text, 3.7% tables, 23.4%

	IMG	TEXT	MM (GME)	MM (T+I)
Avg. Cost (\$)	0.012	0.036	0.022	0.029
Avg. Latency (s)	5.606	7.290	7.897	9.383

Table 6: Average cost of different RAG systems.

1366 figures). Since all domains exhibit similar layout  
1367 patterns, layout does not appear to be a key factor  
1368 in RAG performance.

1369 **F.4 Document page numbers do not affect**  
1370 **performance.**

1371 In the best-performing domains (*commerce man-*  
1372 *ufacturing, education, and legal*), the average  
1373 lengths are 13.1, 14.6, and 12.6 pages, respectively.  
1374 In contrast, the worst-performing domains (*finance,*  
1375 *construction, and healthcare*) average 15.4, 12.9,  
1376 and 12.1 pages. These small differences suggest  
1377 that document length is not a major factor in RAG  
1378 performance.