

# Weekly Report 1

September 6, 2025

## Abstract

This report documents my process of familiarizing myself with the RAG-Anything project through empirical evaluation. I implemented an LLM-as-a-judge framework to assess the system's performance. Initial experiments with LLM-generated QA-pairs from a single PDF, which revealed a need for larger-scale testing. Consequently, I evaluated the system on 230 QA items across text, table, figure, and multimodal types from the MMDocIR benchmark<sup>1</sup>. The results demonstrate the model's relative strength in textual tasks and potential challenges in complex multimodal reasoning, also revealing a need for new evaluation methods for multimodal RAG systems. This work serves as a foundation for understanding the project and guiding future improvements.

## 1 Experiments I

### 1.1 Evaluation Method

I first attempted to evaluate the model with LLM-generated QA-pairs on a single PDF file (the original LightRAG paper).

I modified `process_with_rag()` to extract text content and multimodal items, format context into a QA-pair generation prompt, and call the `llm_model_func()` and `vision_model_func()` from the original code to generate a QA-pair for each text-chunk or multimodal item.

The questions are scored by LLM from 1 to 5 in three dimensions, namely groundness, relevance and stand-alone. The QA-pairs with low-scoring questions are removed.

Then, RAG-Anything is queried with the questions and has the responses collected. The responses and the QA-pair are formatted into an evaluation prompt, utilizing another model to score the RAG response from 1 to 5.<sup>2</sup>

The prompts in this section are adapted from Hugging Face<sup>3</sup>.

### 1.2 Results

Table 1 shows the performance of RAG-Anything on four types of QA-pairs. The "Threshold" in the table refers to the question quality threshold. For instance, when  $\text{Threshold} = 2$ , questions scored below 2 in any of the three dimensions are removed. When  $\text{Threshold} = 1$ , no questions are removed.

The results indicate that RAG-Anything may score higher with better QA-pairs. However, the number of multimodal questions from a single file is too small to be statistically meaningful. Even a minimal increase in the question quality threshold filters out most questions. This indicates that it is difficult to generate high-quality QA-pairs and that there is a need to evaluate the model on larger, pre-reviewed datasets.

---

<sup>1</sup>[https://huggingface.co/datasets/MMDocIR/MMDocIR\\_Evaluation\\_Dataset](https://huggingface.co/datasets/MMDocIR/MMDocIR_Evaluation_Dataset)

<sup>2</sup>Evaluation script: [https://github.com/rnluo/RAG-Anything/blob/main/examples/evaluation\\_on\\_single\\_file.py](https://github.com/rnluo/RAG-Anything/blob/main/examples/evaluation_on_single_file.py)

<sup>3</sup>[https://huggingface.co/learn/cookbook/en/rag\\_evaluation](https://huggingface.co/learn/cookbook/en/rag_evaluation)

Type	Threshold = 1		Threshold = 2	
	Number of QA-pairs	Avg. Score	Number of QA-pairs	Avg. Score
Text	65	4.28	18	<b>4.39</b>
Table	6	4.00	0	N/A
Image	6	3.83	0	N/A
Equation	2	3.00	1	<b>5.00</b>

Table 1: Average LLM-judged score of RAG-Anything on four types of LLM-generated multimodal QA-pairs with two question quality thresholds.

## 2 Experiments II

### 2.1 Evaluation Method

I utilized a subset of the evaluation dataset from MMDocIR(1), a benchmark for multimodal retrieval. The dataset consists of document page screenshots annotated with questions, answers, modality type of the question, and bounding box coordinates of ground truth evidence on pages, etc.

I extracted the first 100 "Pure-text (Plain-text)" questions, 100 "Chart"/"Table"/"Figure" questions, and 30 "multimodal-t"/"multimodal-f" questions. For text questions, I save the provided OCR text into temporary files. For multimodal questions, I crop relevant screenshots with provided bounding box coordinates and save them as temporary files. The temporary files are added into the knowledge base of RAG-Anything, which is queried by corresponding questions. The RAG-responses are scored from 1 to 5 with the same evaluation prompt in section 1.<sup>4</sup>

Since the QA items are based on independent contexts, RAG-Anything must re-initialize and rebuild the knowledge base for every question, which is computationally expensive and time-consuming. This reveals the need for a specially curated dataset in the form of a single, long, multimodal-item-rich document with QA-pairs.

### 2.2 Results

Table 2.2 shows the performance of RAG-Anything on various types of QA items.

Type in original dataset	Number of QA items	Avg. Score
"Pure-text (Plain-text)"	100	3.33
"Chart"	33	2.42
"Table"	26	3.08
"Figure"	41	2.17
"multimodal-t"	25	2.92
"multimodal-f"	5	2.20

Table 2: Average LLM-judged score of RAG-Anything on various types of QA items from the MMDocIR\_Evaluation\_Dataset.

**Pure-text questions** The score of 3.33 is lower than the previous results with LLM-generated QA-pairs. This may result from a discrepancy between the model and the dataset, since RAG-Anything returns long and comprehensive responses and the factual answers in the dataset are much more concise. It reveals a need for datasets specially constructed for RAG systems.

**Multimodal questions** The model scores 3.08 in "Table" questions, significantly higher than in "Chart" and "Figure" questions. This may indicate that the model benefits from the provided table processor.

<sup>4</sup>Evaluation script: [https://github.com/rnluo/RAG-Anything/blob/main/examples/evaluation\\_n\\_on\\_dataset.py](https://github.com/rnluo/RAG-Anything/blob/main/examples/evaluation_n_on_dataset.py)

Table 3: Case Study: One successful case with text and one failed case with multimodal context.

<b>Type</b>	"Pure-text (Plain-text)"								
<b>Context</b>	[OCR text...]								
<b>Question</b>	What was the value of absolute percentage shortfall in India's GDP growth from 2002 to 2003 compared to the forecast?								
<b>Answer</b>	2.4%								
<b>RAG-Anything</b>	<p>The absolute percentage shortfall in India's GDP growth from 2002 to 2003 compared to the forecast can be calculated as follows:</p> <ul style="list-style-type: none"> <li>- Forecast GDP growth: 6.7%</li> <li>- Actual GDP growth: 4.3%</li> </ul> <p>The absolute percentage shortfall is calculated as:</p> $\text{Shortfall} = \text{Forecast} - \text{Actual} = 6.7\% - 4.3\% = 2.4\%$ <p>Thus, the absolute percentage shortfall in India's GDP growth was <b>2.4%</b> from 2002 to 2003.</p> <p><b>References</b> [DC] tmpji3zt16n.txt</p>								
<b>Type</b>	"multimodal-t"								
<b>Context</b>	<table border="1"> <thead> <tr> <th>Setting</th><th>F1</th></tr> </thead> <tbody> <tr> <td>Distractor</td><td>67.08</td></tr> <tr> <td>Open-domain 10 Paragraphs</td><td>38.40</td></tr> <tr> <td>Open-domain 500 Paragraphs + Gold Paragraph</td><td>39.12</td></tr> </tbody> </table>	Setting	F1	Distractor	67.08	Open-domain 10 Paragraphs	38.40	Open-domain 500 Paragraphs + Gold Paragraph	39.12
Setting	F1								
Distractor	67.08								
Open-domain 10 Paragraphs	38.40								
Open-domain 500 Paragraphs + Gold Paragraph	39.12								
<b>Question</b>	By how much does the F1 score improve when the gold paragraph is added to the open-domain setting with 500 paragraphs?								
<b>Answer</b>	The F1 score improves by 14.00 when the gold paragraph is added.								
<b>RAG-Anything</b>	<p><b>F1 Score Improvement Overview</b></p> <p>Based on the information available, there are various numerical values provided in the context. The relevant values that seem to relate to performance measurements, possibly including the F1 score, are 38.40, 39.12, 53.12, and 67.08. These numbers might represent metrics used in evaluating an open-domain text composition task, such as the ones mentioned regarding different paragraph settings [KG] Open-domain 10 Paragraphs - Open-domain 500 Paragraphs + Gold Paragraph.</p> <p>Although the exact nature of the improvement in F1 score is not explicitly stated, we can infer from the context provided that an enhancement might be associated with the addition of the 'Gold Paragraph' in the open-domain settings involving 500 paragraphs. The exact degree of improvement, however, would require detailed data that is not present in the provided context.</p> <p><b>Conclusion</b></p> <p>It is not possible to definitively determine by how much the F1 score improves when the gold paragraph is added based on the information available as it lacks specific comparative details directly related to the F1 score.</p> <p><b>References</b></p> <ul style="list-style-type: none"> <li>- [KG] Open-domain 10 Paragraphs - Open-domain 500 Paragraphs + Gold Paragraph</li> <li>- [KG] 38.40 - 53.12</li> <li>- [KG] 39.12 - 67.08</li> <li>- [DC] tmpx9v9v4oi.png</li> </ul>								

There is another reason why the model performs worse on multimodal questions than on pure-text questions, illustrated in Table 3. The two cases both involve subtracting two figures to obtain the answer. While the model succeeds with clear reasoning in the pure-text case, it fails to comprehend the task in the multimodal case, even though the figures are correctly extracted. This may indicate that the model can be improved in complex multimodal reasoning, especially in the effective comprehension of short texts in multimodal items.

## References

- [1] Kuicai Dong, Yujing Chang, Xin Deik Goh, Dexun Li, Ruiming Tang, and Yong Liu. Mmdocir: Benchmarking multi-modal retrieval for long documents, 2025.