

# NeuSym-RAG: Hybrid Neural Symbolic Retrieval with Multiview Structuring for PDF Question Answering

Anonymous ACL submission

## Abstract

The increasing number of academic papers poses significant challenges for researchers to efficiently acquire key details. While retrieval augmented generation (RAG) shows great promise in large language model (LLM) based automated question answering, previous works often isolate neural and symbolic retrieval despite their complementary strengths. Moreover, conventional single-view chunking neglects the rich structure and layout of PDFs, e.g., sections and tables. In this work, we propose NeuSym-RAG, a hybrid neural symbolic retrieval framework which combines both paradigms in an interactive process. By leveraging multi-view chunking and schema-based parsing, NeuSym-RAG organizes semi-structured PDF content into both the relational database and vectorstore, enabling LLM agents to iteratively gather context until sufficient to generate answers. Experiments on three full PDF-based QA datasets, including a self-annotated one AIRQA-REAL, show that NeuSym-RAG stably defeats both the vector-based RAG and various structured baselines, highlighting its capacity to unify both retrieval schemes and utilize multiple views.

## 1 Introduction

With the exponential growth in academic papers, large language model (LLM) based question answering (QA) systems show great potential to help researchers extract key details from emerging studies. However, individual PDFs often exceed prompt limits, and user queries may span multiple documents. To tackle these challenges, retrieval-augmented generation (RAG, Huang and Huang, 2024) is effective for knowledge-intensive QA.

Despite its wide application, the classic *neural retrieval* (Guu et al., 2020) often fails when handling precise queries involving mathematical operations, comparisons, or aggregations. For example, in the top-left of Figure 1, the total number of tables cannot be determined through retrieved chunks

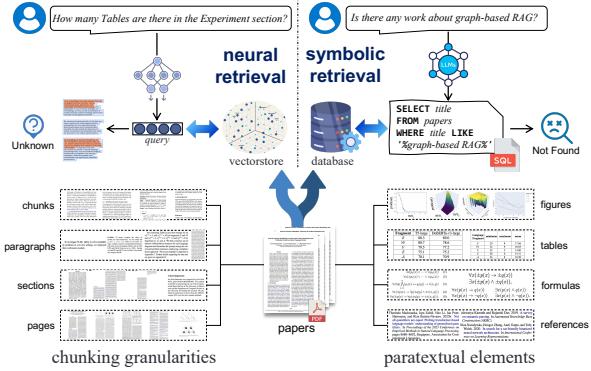


Figure 1: Motivation of the proposed NeuSym-RAG. (Upper) Two paradigms of retrieval strategies. (Bottom) PDF documents can be split based on different granularities and they contain many paratextual elements.

as they are scattered across the document. On the other hand, *symbolic retrieval* such as TAG (Biswal et al., 2024) relies on semantic parsing (Berant et al., 2013) techniques, e.g., text-to-SQL (Li et al., 2024b), to directly extract the target information from the structured database. Unfortunately, such precise queries often break down in semantic fuzzy matching or morphological variations, e.g., “*graph-based RAG*” versus “*GraphRAG*”. Previous literature investigates these two paradigms in isolation.

Furthermore, the most widely utilized scheme to segment documents into chunks is based on a fixed length of consecutive tokens, possibly considering sentence boundaries or more delicate granularities (Yang, 2023). However, for semi-structured PDF documents of research papers, this common practice neglects the intrinsic structure of sections and the salient features of paratextual tables and figures (illustrated at the bottom of Figure 1). The distinct layout of PDF files offers a more structured view towards segmenting and arranging content.

To this end, we propose a hybrid **Neu**ral **Sy**mbolic retrieval framework (NeuSym-RAG) for PDF-based question answering, which combines both retrieval paradigms into an interactive pro-

cedure. During pre-processing, the PDF file of each paper is fed into a pipeline of parsing functions (§ 2.2) to complete the metadata, segment raw texts based on different views, and extract various embedded paratextual elements (e.g., tables and figures). These identified elements are populated into a relational database, specifically designed for semi-structured PDF. Then, we select encodable column values from the populated database for storage into the vectorstore (§ 2.3). These cell values present diverse views in interpreting PDF content. Moreover, the database schema graph is leveraged to organize the vectors into a well-formed structure. To answer the user question, LLM agents can adaptively predict executable actions (§ 2.4.1) to retrieve desired information from either backend environment (database or vectorstore) in a multi-round fashion. This agentic retrieval terminates when the collected information suffices to answer the input question, in which case LLM agent will predict a terminal “GENERATEANSWER” action.

To validate NeuSym-RAG, we convert and annotate three full PDF-based academic QA datasets, including a self-curated one AIRQA-REAL with 553 samples and 18 instance-specific evaluation metrics (§ 3.1). Experiments on both closed- and open-source LLMs demonstrate that it remarkably outperforms the classic neural RAG (17.3% points on AIRQA-REAL) and various structured methods like HybridRAG (Sarmah et al., 2024) and GraphRAG (Edge et al., 2024). The ablation study highlights: 1) the positive impact of model size increase on agentic retrieval, 2) the superiority of integrating multiple chunking perspectives, and 3) the low sensitivity to LLM choice for pre-processing and high sensitivity for subjective evaluation.

To sum up, our contributions are three-fold:

- We are the first to integrate both vector-based neural retrieval and SQL-based symbolic retrieval into a unified and interactive NeuSym-RAG framework through executable actions.
- We incorporate multiple views for parsing and vectorizing PDF documents, and adopt a structured database schema to systematically organize both text tokens and encoded vectors.
- Experiments on three realistic full PDF-based QA datasets *w.r.t.* academic research validate the superiority of NeuSym-RAG over various neural and symbolic baselines.

## 2 Framework of NeuSym-RAG

This section presents the complete framework of NeuSym-RAG, comprising three stages (Figure 2).

### 2.1 The Overall Procedure

To answer user questions about a semi-structured PDF file, the entire workflow proceeds as follows:

- 1) **Parsing:** Firstly, we pass the raw PDF file into a pipeline of functions to segment it in multi-view, extract non-textual elements, and store them in a schema-constrained database (DB).
- 2) **Encoding:** Next, we identify those encodable columns in the DB, and utilize embedding models for different modalities to obtain and insert vectors of cell values into the vectorstore (VS).
- 3) **Interaction:** Finally, we build an iterative Q&A agent which can predict executable actions to retrieve context from the backend environment (either DB or VS) and answer the input question.

### 2.2 Multiview Document Parsing

At this stage, for each incoming PDF file, we aim to parse it with different perspectives into a relational database DuckDB (Mühleisen and Raasveldt, 2024). The pipeline of parsing functions includes (in the top middle of Figure 2): 1) Querying scholar APIs (e.g., arxiv) to obtain the metadata such as the authors and published conference, such that we can support metadata-based filtering (Poliakov and Shvai, 2024) during retrieval. 2) Splitting the text based on different granularities with tool PyMuPDF (Artifex Software, 2023), e.g., pages, sections, and fixed-length continuous tokens. 3) Leveraging OCR models to extract non-textual elements (we choose the tool MinerU, Wang et al., 2024a). 4) Asking large language models (LLMs) or vision language models (VLMs) to generate concise summaries of the parsed texts, tables, and images (Xu et al., 2023). The retrieved metadata, parsed elements and predicted summaries will all be populated into the symbolic DB. We handcraft the DB schema in advance, which is carefully designed and universal for PDF documents (see the middle part of Figure 2 or App. E for visualization).

### 2.3 Multimodal Vector Encoding

After the completion of the symbolic part, we switch to the neural encoding of parsed elements. Firstly, we label each column in the DB schema

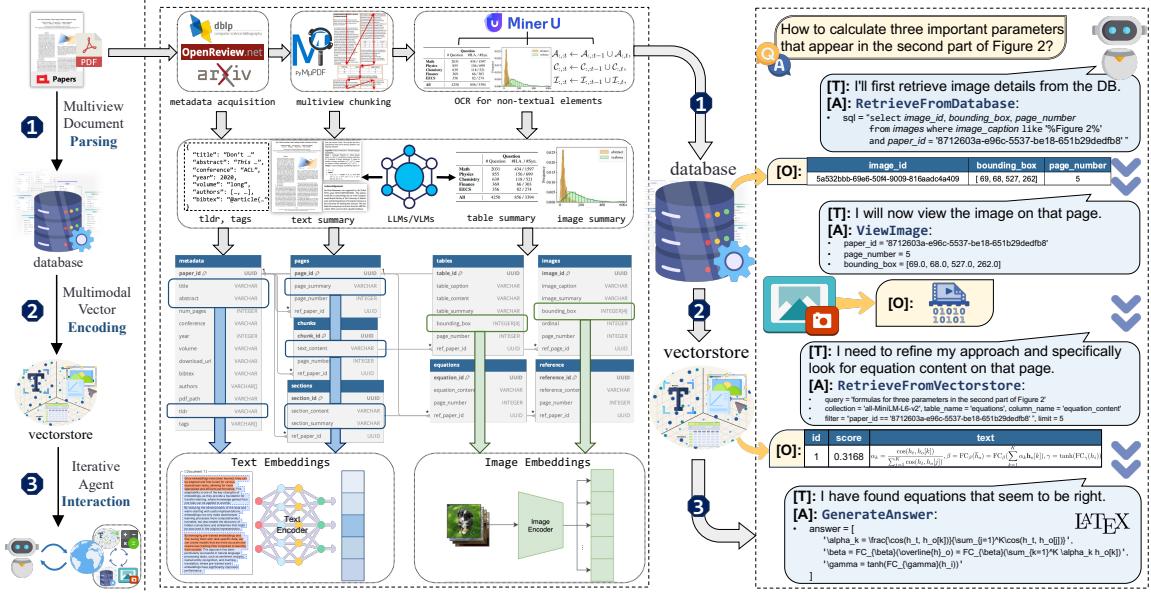


Figure 2: Overview of the NeuSym-RAG framework. The demonstration example comes from a real and simplified case in our labeled dataset (§ 3.1). “[T]”, “[A]”, and “[O]” represent thought, action and observation, respectively.

as “encodable” or not. For text modality, we mark those columns of data type “`varchar`” as encodable whose cell values are relatively long (e.g., column “`pages.page_summary`” in Figure 3); while for images, the “`bounding_box`” columns which record 4 coordinates of figures or tables are recognized to extract the square region in the PDF. Next, we resort to encoding models of both modalities to vectorize text snippets or cropped images.

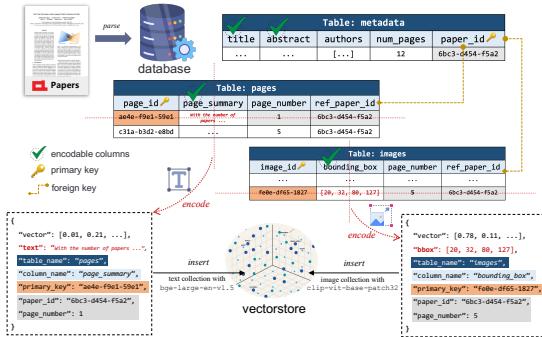


Figure 3: Illustration of how to convert an encodable cell value in the database to one data entry in the vectorstore.

To build a one-to-one mapping between each cell value in the DB and neural vector in the VS, we supplement each data point in the VS with its corresponding table name, column name, and primary key value for that row in the DB. This triplet can uniquely identify each value in the DB. Besides, we add 2 extra fields, namely “`paper_id`” and “`page_number`”, into the JSON dict to enable metadata filtering. These data entries will be in-

serted into the VS, categorized into different collections based on the encoding model and modality<sup>1</sup>.

Through the first two stages (§ 2.2 and § 2.3), various chunking perspectives in the VS are intrinsically connected via a structured DB schema (visualized in Figure 10). Moreover, long-form texts or visual bounding boxes in the DB are vectorized into the VS to support fuzzy semantic matching.

## 2.4 Iterative Agent Interaction

Now that the database and vectorstore have been populated, we can build an RAG agent which proactively retrieves context from both the DB and VS.

### 2.4.1 Action Space Design

Firstly, we introduce the 5 parameterized actions with arguments that agents can take during interaction, namely `RETRIEVEFROMVECTORSTORE`, `RETRIEVEFROMDATABASE`, `VIEWIMAGE`, `CALCULATEEXPR`, and `GENERATEANSWER`.

**RETRIEVEFROMVECTORSTORE** This action converts the classic static retrieval into real-time dynamic retrieval. It has multiple arguments which are adjustable (see Lst. 1), e.g., “`query`” encourages LLMs to rewrite the user intention more clearly (Ma et al., 2023), while “`table_name`” and

<sup>1</sup>We create 4 collections in the Milvus (Wang et al., 2021) vectorstore: for text embeddings, we use BM25 (Robertson and Zaragoza, 2009), all-MiniLM-L6-v2 (Wang et al., 2020), and bge-large-en-v1.5 (Xiao et al., 2023), while clip-vit-base-patch32 (Radford et al., 2021) for images.

“column\_name” request the agent to select an appropriate perspective for retrieval. During execution, the environment will retrieve context concerning “query” with specified constraints “filter” (e.g., page\_number=1) from the VS and return the observation in a tabular format (see App. F.2.3).

```

1  RetrieveFromVectorstore(
2      # user input can be rephrased
3      query: str,
4      # select encoding model/modality
5      collection_name: str,
6      # (table_name, column_name) together
7          # defines which view to search
8      table_name: str,
9      column_name: str,
10     # allow fine-grained meta filtering
11     filter: str = '',
12     limit: int = 5
)

```

Listing 1: RETRIEVEFROMVECTORSTORE action with its parameters in function calling format.

**RETRIEVEFROMDATABASE** This action accepts a single parameter “sql”. It operates by executing the provided SQL query against the prepared DB and fetching the resulting table to the agent. Similar to ToolSQL (Wang et al., 2024b), through iterative symbolic retrieval, agents can predict different SQLs to explore the DB and exploit multiple pre-parsed views on interpreting PDF content.

**VIEWIMAGE** To address potential PDF parsing errors during the first stage (§ 2.2) and leverage features of more advanced VLMs, we devise this action to extract a cropped region (defined by “bounding\_box”) in one PDF page and send the base64 encoded image back to agents. Notably, the concrete coordinates can be obtained through retrieval from either the DB or VS during interaction. This way, the agent will acquire the desired image as observation and then reason based on it.

```

1  ViewImage(
2      paper_id: str,
3      page_number: int,
4      # 4-tuple of float numbers, if [], 
5          # return the image of entire page
6      bounding_box: List[float] = []
)

```

Listing 2: VIEWIMAGE action with its parameters in function calling format.

**CALCULATEEXPR** Inspired by Xu et al. (2024), we also include this action which accepts a Python “expr” (e.g.,  $2 + 3 * 4$ ) and returns the calculation result. This simple integration stably reduces hallucination w.r.t. math problems in our pilot study.

**GENERATEANSWER** This is the *terminal* action that returns the final answer when agents determine that the retrieved context suffices to resolve the input question. The only parameter “answer” can be of any type depending on the user requirement. Refer to App. F.2.1 for formal action definitions.

## 2.4.2 Hybrid Neural Symbolic Retrieval

In each turn, agents predict one action (see App. F.2.2 for different action formats) to interact with the environment and obtain the real-time observation. We adopt the popular ReAct (Yao et al., 2023) framework, which requires agents to output their thought process first (see the right part of Figure 2). This iterative retrieval (Shao et al., 2023) enables agents to leverage complementary strengths of both retrieval paradigms in continuous steps.

For example, agents can firstly filter relevant rows by executing a SQL program in the database. RETRIEVEFROMDATABASE is particularly adept at handling structured queries and metadata-based constraints. After extracting primary key values of rows for the intermediate output, we can further narrow down the final result set by inserting those key values into the “filter” parameter of action RETRIEVEFROMVECTORSTORE to conduct neural semantic matching, i.e.,

`filter='primary_key in [pk values of sql exec]'.`

Conversely, agents can firstly query the VS via neural search to select the most relevant entries, and then treat this transitional set as a temporary table or condition in the subsequent SQL retrieval.

## 3 Experiment

This section introduces our human-labeled dataset AIRQA-REAL, the main experiment and ablation study. Code and data are publicly available <sup>2</sup>.

### 3.1 Q&A Dataset on AI Research Papers

**AIRQA-REAL** Previous Q&A datasets on academic papers typically focus on simple questions based on a single page or merely the abstract of the PDF, which is far from reflecting real-world scenarios. To address this, we annotate a more complex Q&A dataset regarding full and even multiple PDFs, called AIRQA-REAL. Based on published AI papers in 2023 and 2024, 16 researchers manually annotate 553 questions which span across: a) 3 task types (i.e., single-doc details, multi-doc

<sup>2</sup><https://anonymous.4open.science/r/neu-sym-rag/>

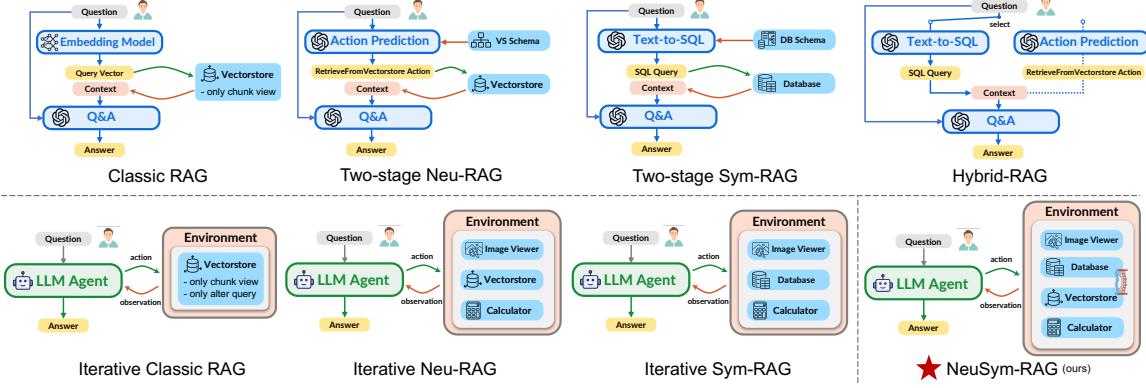


Figure 4: Comparison of different structured baselines with our NeuSym-RAG. Refer to App. C for text description.

Category	Question
text	What are the main components of ERRA model?
table	On the dataset proposed by this work, how much does the GPT-3.5-turbo model improve its GPT4score after using Graph-CoT?
image	Considering the performance of ChatDev agent on DSEval-LeetCode benchmark, what is the most common cause of the errors?
formula	How does Multi-DYLE combine the three different losses as the objective of training?
metadata	Which conference was the paper 'Fact-Checking Complex Claims with Program-Guided Reasoning' published in? Is it a long, short or findings?

Table 1: Examples of single-doc type from AIRQA-REAL dataset. See App.A.2 for the other two types.

analysis, and paper retrieval), b) 5 categories (i.e., text, table, image, formula, and metadata), and c) 2 evaluation genres (i.e., hard-coding objective metrics and LLM-based subjective assessment). We showcase one example for each category in Table 1.

**Other Benchmarks** To validate the universality, we also convert two other full-PDF-based Q&A datasets, namely M3SciQA (Li et al., 2024a) and SciDQA (Singh et al., 2024). These two benchmarks have 452 and 2937 test samples, respectively. And we adjust them into the data format of AIRQA-REAL. See App. A for more details about datasets.

**Evaluation Metrics** While most long-form Q&A datasets utilize LLMs or human experts for accuracy (including the aforementioned M3SciQA and SciDQA), we propose more flexible and precise evaluation metrics. Specifically, we 1) adopt answer formatting similar to DABench (Hu et al., 2024), which imposes Python-style restrictions on output (e.g., “Your answer should be a Python list of two strings.”), and 2) implement instance-specific, execution-based evaluation by designing

18 functions with optional parameters. These functions are categorized as either subjective or objective, depending on whether they involve the judgement from LLMs (see App. D).

### 3.2 Experiment Settings

**Baselines** We implement 11 methods (partly illustrated in Figure 4), ranging from trivial input questions, to more advanced structured models. The main differences lie in: a) Single-view or Multi-view: the Classic RAG only has access to text chunks, while the others can select from multiple pre-parsed perspectives; b) Two-stage or Iterative: the major distinction between the two rows in Figure 4 is whether agentic retrieval is permitted; c) Composition of the environment: the complete backend contains relational database, vectorstore, image viewer and calculator. Different baselines may only include part of them, which further restricts the permissible actions (see App. F.2.1).

**LLMs and Hyper-parameters** We evaluate NeuSym-RAG with various LLMs. For closed-source ones, we use GPT-4o-mini-2024-07-18 and GPT-4-1106-vision-preview. Regarding state-of-the-art open-source LLMs, we include InternLM2.5 (Cai et al., 2024), Llama-3.3-70B-Instruct (Hugging Face Team, 2023), Qwen2.5-VL-72B-Instruct (Team, 2025), and DeepSeek-R1 (Guo et al., 2025). As for hyper-parameters, the temperature is set to 0.7 and top\_p is fixed to 0.95. The maximum retrieved tokens in each turn and the cutoff for full-text input are both limited to 5k. And the threshold of interaction turns is 20.

### 3.3 Main Experiment

In Table 2, we show performances of our proposed method on different LLMs and datasets. Accord-

304  
305  
306  
307  
308

309  
310  
311  
312  
313  
314  
315

316  
317  
318  
319  
320  
321  
322  
323  
324  
325

326  
327  
328  
329  
  
330  
331  
332  
333  
334  
335  
  
336  
337  
338  
339  
  
340  
341  
342  
343  
344  
  
345  
346  
347  
348  
349  
350  
  
351  
352  
353  
354  
355  
356  
357  
  
358  
359  
360

Model	AIRQA-REAL						M3SciQA			SciDQA			
	text	table	image	formula	metadata	Avg	table	image	Avg	table	image	formula	Avg
Classic-RAG													
GPT-4o-mini	12.3	11.9	12.5	16.7	13.6	13.4	17.9	10.6	15.6	59.4	60.4	59.3	59.8
GPT-4V	13.2	13.9	10.0	13.9	13.6	14.7	12.1	8.8	11.1	56.6	56.8	58.1	57.4
Llama-3.3-70B-Instruct	8.7	7.9	9.5	16.7	0.0	10.0	12.7	8.1	11.3	56.8	58.8	58.9	58.0
Qwen2.5-VL-72B-Instruct	9.6	5.9	11.9	11.1	13.6	10.5	11.6	11.6	11.6	54.8	56.9	56.3	56.2
DeepSeek-R1	11.7	13.9	9.5	30.6	9.1	13.9	11.9	9.5	11.2	63.9	61.3	61.7	62.4
NeuSym-RAG													
GPT-4o-mini	33.0	12.9	11.9	19.4	18.2	30.7	18.7	16.6	18.0	63.0	63.6	62.5	63.0
GPT-4V	38.9	<b>18.8</b>	<b>23.8</b>	<b>38.9</b>	<b>27.3</b>	37.3	13.7	13.4	13.6	62.6	63.5	63.2	63.1
Llama-3.3-70B-Instruct	30.6	11.9	16.7	16.7	<b>27.3</b>	29.3	<b>26.3</b>	17.6	<b>23.6</b>	55.5	57.3	56.6	56.4
Qwen2.5-VL-72B-Instruct	<b>43.4</b>	15.8	11.9	25.0	<b>27.3</b>	<b>39.6</b>	20.2	<b>22.7</b>	21.1	60.2	60.6	61.8	60.5
DeepSeek-R1	33.2	16.8	11.9	27.8	18.2	32.4	19.0	13.7	17.4	<b>64.3</b>	<b>64.6</b>	<b>63.9</b>	<b>64.5</b>

Table 2: Main results of different LLMs using Classic-RAG or NeuSym-RAG on three datasets.

Method	Neural	Symbolic	Multi-view	# Interaction(s)	sgl.	multi.	retr.	subj.	obj.	Avg
Question only				1	5.7	8.0	0.4	9.4	2.7	4.0
Title + Abstract	✗	✗	✗	1	5.7	14.0	0.0	13.1	3.6	5.4
Full-text w/. cutoff				1	28.3	10.7	0.4	26.2	7.6	11.2
Classic RAG	✓	✗	✗	1	18.2	4.0	9.4	8.4	11.0	10.5
Iterative Classic RAG	✓	✗	✗	≥ 2	8.2	10.0	15.2	5.6	13.2	11.8
Two-stage Neu-RAG	✓	✗	✓	2	19.5	10.0	5.3	15.9	9.4	10.7
Iterative Neu-RAG	✓	✗	✓	≥ 2	<b>37.7</b>	18.7	48.4	<b>32.7</b>	38.3	37.3
Two-stage Sym-RAG	✗	✓	✓	2	12.2	5.4	9.4	10.6	8.7	9.1
Iterative Sym-RAG	✗	✓	✓	≥ 2	32.1	14.7	33.6	27.1	28.3	28.0
Graph-RAG	✓	✗	✓	2	22.2	11.1	0.0	21.1	11.5	15.6
Hybrid-RAG	✓	✓	✓	2	23.3	9.3	5.7	16.8	10.5	11.8
<b>NeuSym-RAG (ours)</b>	✓	✓	✓	≥ 2	28.3	<b>32.3</b>	<b>58.2</b>	27.1	<b>42.6</b>	<b>39.6</b>

Table 3: Performances of different RAG methods on AIRQA-REAL dataset, where “# Interaction(s)” denotes the number of LLM calls for a single question. See App. C for detailed introduction of each RAG baseline.

1) NeuSym-RAG remarkably outperforms the Classic RAG baseline across all datasets, with a minimal 17.3% improvement on AIRQA-REAL for all LLMs. With more customizable actions and more trials, NeuSym-RAG can interact with the backend environment with higher tolerance and learn through observations to retrieve proper information for question answering. 2) VLMs perform better in tasks that require vision capability, e.g., in M3SciQA where LLMs have to view an anchor image in the first place. 3) Open-source LLMs are capable of handling this complicated interactive procedure in a zero-shot paradigm, and even better than closed-source LLMs. NeuSym-RAG with Qwen2.5-VL-72B-Instruct elevates the overall accuracy by 29.1%, compared to the Classic RAG baseline. Surprisingly, it surpasses GPT-4o-mini by an impressive 8.9%, which highlights the universality of NeuSym-RAG.

Next, to figure out the contribution of each component in NeuSym-RAG, we compare the performances of different structured RAG agent methods described in § 3.2. From Table 3, we can observe that: 1) Two-stage Neu-RAG outperforms Clas-

sic RAG, while Hybrid RAG achieves even further improvements. This consistent increase can be attributed to the fact that agents can adaptively determine the parameters of actions, e.g., which chunking view to select for neural retrieval. 2) Iterative retrievals are superior to their two-stage variants. Through multi-turn interaction, the agent can explore the backend environment and select the most relevant information to answer the question. 3) As the number of interactions increases, objective scores rise faster than subjective scores, indicating that with more retrievals, LLMs generate more rationally. On the whole, NeuSym-RAG defeats all adversaries. It verifies that providing multiple views and combining two retrieval strategies both contribute to the eventual performance.

### 3.4 Ablation Study

For brevity, in this section, GPT denotes GPT-4o-mini and Qwen represents Qwen2.5-72B-Instruct. And unless otherwise specified, we utilize Qwen.

**Different Model Scales** In Figure 5, we show the performance changes by varying the model scales. Accordingly, we observe that: 1) The per-

361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379

385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407

408 performances consistently improve as the scales of  
 409 LLMs increase. 2) Unfortunately, the R1-Distill-  
 410 Qwen series models lag far behind their counter-  
 411 parts (Qwen2.5-Instruct). Although they generate  
 412 valuable thoughts during the interaction, they strug-  
 413 gle to exactly follow the action format specified  
 414 in our prompt. We hypothesize that they severely  
 415 overfit the instructions during R1 post-training.

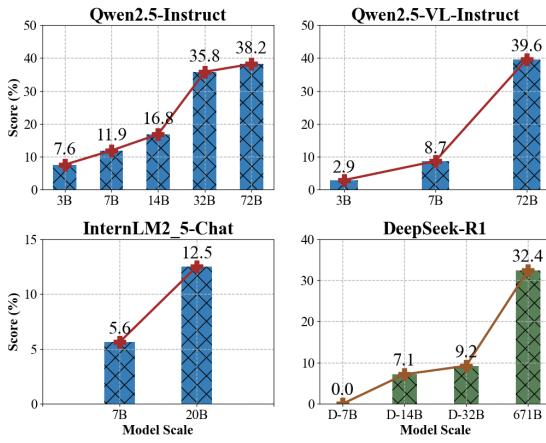


Figure 5: Performances of different model sizes on AIRQA-REAL dataset. “D-” denotes “Distill-Qwen-”.

416 **Chunking with Multiple Views** In the Classic  
 417 RAG baseline, we segment raw documents follow-  
 418 ing continuous 512 tokens. We wonder whether  
 419 different chunking views lead to varied outcomes.  
 420 Results in Table 4 are insightful: the classic chunk-  
 421 ing strategy is still the best in general. However,  
 422 chunking views tailored to specific aspects may  
 423 achieve the best results in their particular categories.

Retrieved Column	text	table	image	formula	Overall (%)
chunk_content	13.95	<b>15.00</b>	<b>22.22</b>	28.57	<b>18</b>
section_content	<b>18.60</b>	10.00	5.56	14.29	12
table_content	13.95	<b>15.00</b>	0.00	9.52	9
image_summary	4.65	0.00	16.67	4.76	6
equation_content	6.98	0.00	5.56	<b>33.33</b>	8

Table 4: Performances of different retrieval options of the Classic RAG baseline on a subset of AIRQA-REAL.

Agent Model	Summary Model	objective	subjective	Overall (%)
GPT	GPT	24.62	34.29	28
	Qwen	23.08	34.29	27
Qwen	GPT	30.77	57.14	40
	Qwen	30.77	54.28	39

Table 5: Performances of different models with different LLM summary models on a subset of AIRQA-REAL.

**Details of Pre-processing and Evaluation** Aiming to reduce the high cost of API calls, we explore the impact of using open-source v.s. closed-source LLMs during 1) PDF pre-parsing (the summary generation), and 2) subjective evaluation. As shown in Table 5, cheaper open-source LLMs appear to be a viable alternative to closed-source ones, since the final results are insensitive to element summaries produced by different LLMs. As for evaluation, we randomly sample 100 long-form answers requiring subjective assessment and recruit AI researchers to independently measure the semantic equivalence compared to the reference answer. Figure 6 indicates that, in most cases, LLM evaluations align closely with human judgments. Notably, closed-source GPT-4o aligns exactly with human judgments, whereas open-source LLMs are occasionally either too strict or too lenient. Therefore, we choose closed-source LLMs for subjective evaluation throughout the experiments.

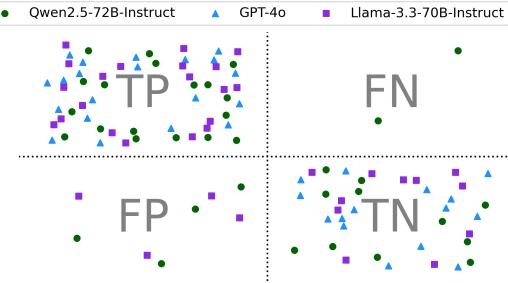


Figure 6: The consistency between human and LLM-based evaluation. We classify examples that are judged as correct by human as positive samples.

### 3.5 Case Study

In Table 6, we serialize the interaction trajectory of 3 iterative methods for the same question. NeuSym-RAG firstly predicts a RETRIEVEFROMVECTORSTORE action to search chunks regarding “WER values”. After detecting relevant figures, it predicts a RETRIEVEFROMDATABASE action to access their captions and summaries. In contrast, Iterative Neu-RAG fails to exploit the figures due to the absence of the DB, and Iterative Sym-RAG struggles to locate relevant context. In this case, we can clearly observe one benefit or mode of combining both retrieval actions: the VS helps agents identify relevant context, while the DB for precise queries. This synergy achieves the accurate answer unattainable by either retrieval type alone, emphasizing the significance of their collaboration.

Iterative Neu-RAG	
Turn #1	
[Thought]:	I will now try to retrieve the relevant sections by focusing on the content that mentions WER values for DD2 and NV1.
[Action]:	RETRIEVEFROMVECTORSTORE (query="WER values for DD2 and NV1 methods", table_name="sections", column_name="section_content", filter="pdf_id == '81b6a6b0-a195-5cae-9e30-137150b64352'")
[Observation]:	{"text": "Figures 3 and 4 depict WER and PER values respectively with ASR on the six different methods. As expected, artificial speech with variations obtained variation, thus confirming that Wav2Vec 2.0 peripher WER scores ≈ +0.57 than speech without forms better on speech ..."} -----
Turn #2	
[Thought]:	The difference in WER values between DD2 and NV1 is 0.57.
[Action]:	GENERATEANSWER (answer=0.57) ✗
Iterative Sym-RAG	
Turn #1	
[Thought]:	I will now check the 'sections' table to see if the WER values are mentioned in any of the sections of the paper.
[Action]:	RETRIEVEFROMDATABASE (sql="SELECT section_title, section_content FROM sections WHERE ref_paper_id = '81b6a6b0-a195-5cae-9e30-137150b64352';")
[Observation]:	{"section_title": "5 Experiments", "section_content": "The experiments aim to evaluate the extent to which it is possible to infer accent-related phonemes confusions from artificially accented speech. For that purpose, ..."} -----
Turn #2	
[Thought]:	I will now try to extract the WER values from the '5 Experiments' section by parsing the text content.
[Action]:	RETRIEVEFROMDATABASE (sql="SELECT section_content FROM sections WHERE ref_paper_id = '81b6a6b0-a195-5cae-9e30-137150b64352' AND section_title = '5 Experiments';")
[Observation]:	{"section_content": "The experiments aim to evaluate the extent to which it is possible to infer accent-related phonemes confusions from artificially accented speech. For that purpose, ..."} -----
Turn #3	
[Thought]:	I will assume the WER values are as follows: - DD2: 25.0 - NV1: 10.0
	The difference in WER values is 15.0.
[Action]:	GENERATEANSWER (answer=15.0) ✗
NeuSym-RAG	
Turn #1	
[Thought]:	I will perform a similarity search in the vectorstore to find the relevant context that contains the WER values.
[Action]:	RETRIEVEFROMVECTORSTORE (query="WER values for DD2 and NV1 in the paper with id 81b6a6b0-a195-5cae-9e30-137150b64352", table_name="sections", column_name="section_content", filter="pdf_id == '81b6a6b0-a195-5cae-9e30-137150b64352'")
[Observation]:	{"text": "Figures 3 and 4 depict WER and PER values respectively with ASR on the six different methods. As expected, ..."} -----
Turn #2	
[Thought]:	I will now query the 'figures' table to find the figure captions that might contain the WER values.
[Action]:	RETRIEVEFROMDATABASE (sql="SELECT image_caption FROM images WHERE ref_paper_id = '81b6a6b0-a195-5cae-9e30-137150b64352' AND (image_caption LIKE '%Figure 3%' OR image_caption LIKE '%Figure 4%');")
[Observation]:	{"image_caption": "Figure 3: WER scores for DD2 vs other methods.", "image_caption": "Figure 4: PER scores for DD2 vs other methods."} -----
Turn #3	
[Thought]:	I will try to retrieve the detailed content of the figures from the 'images' table in the database.
[Action]:	RETRIEVEFROMDATABASE (sql="SELECT image_summary FROM images WHERE ref_paper_id = '81b6a6b0-a195-5cae-9e30-137150b64352' AND image_caption LIKE '%Figure 3%'")
[Observation]:	{"image_summary": "The bar chart compares Word Error Rates across six methods: NV1 (0.131), NV2 (0.077), DD1 (0.741), KB1 (0.658), KB2 (0.666), and DD2 (0.729). NV2 has the lowest error rate, while DD1 has the highest."} -----
Turn #4	
[Thought]:	I will extract the WER values for DD2 and NV1 and calculate the difference.
[Action]:	CALCULATEEXPR (expr="0.729 - 0.131")
[Observation]:	The calculated result is: 0.598 -----
Turn #5	
[Thought]:	The difference in WER values between DD2 and NV1 is 0.598.
[Action]:	GENERATEANSWER (answer=0.598) ✓

Table 6: Case study on AIRQA-REAL dataset. The question of the case is “In terms of WER values with ASR across the six different methods tested in the paper, how much higher is DD2 compared to NV1?”.

## 4 Related Work

**Structured and Iterative Retrieval** Previous literature has proposed various RAG methods for structured knowledge sources. GraphRAG (Edge et al., 2024) builds a knowledge graph (KG) to summarize partial and final answers. HybridRAG (Sarmah et al., 2024) integrates KG-based and vector-based methods for financial documents. FastRAG (Abane et al., 2024) employs schema and script learning to extract structured data, combining text search with KG. Another promising line of work proves that agents can benefit from iterative retrieval. IRCoT (Trivedi et al., 2023) handles multi-step QA by interleaving retrieval with chain-of-thought. Iter-RetGen (Shao et al., 2023) integrates retrieval and generation iteratively for collective knowledge. CoA (Zhang et al., 2024) uses multiple agents for long-context tasks, with a manager agent combining their results. Our NeuSym-RAG integrates neural and symbolic retrieval in an agentic interaction to fetch multi-view context.

**PDF-based Q&A Benchmarks** Existing question answering datasets over PDF documents often overlook the underlying structure and layout. QASPER (Dasigi et al., 2021) generate questions from merely titles and abstracts, while FinQA (Chen et al., 2021) targets single-page Q&A. QASA (Lee et al., 2023) retains section structures but misses other key elements like figures and tables. More recent works, like Visconde (Pereira et al., 2023), M3SciQA (Li et al., 2024a), and SciDQA (Singh et al., 2024), focus on long-form subjective evaluation. This work addresses these limitations by tackling questions from full PDFs of AI research papers, including objective metrics.

## 5 Conclusion

In this work, we propose a hybrid neural symbolic retrieval framework (NeuSym-RAG) for semi-structured PDF-based question-answering (Q&A). Experiments on three human-labeled realistic datasets regarding AI research, especially on the self-annotated AIRQA-REAL, show that NeuSym-RAG remarkably outperforms the Classic RAG baseline by more than 17.3% points. Future works include: 1) training a specialized LLM agent of medium size to further improve performances and efficiency, and 2) extending Q&A tasks to other vertical domains such as healthcare and law which heavily rely on external well-structured PDF files.

## 511 Limitations

512 Despite the superiority of the proposed NeuSym-  
513 RAG framework, it still has the following limita-  
514 tions: 1) The pipeline of parsing and encoding the  
515 PDF content is relatively slow, especially when  
516 higher parsing accuracy is demanded. That is, we  
517 need OCR models typically with larger capacity  
518 and longer processing time. This method is suitable  
519 if all documents are pre-stored in a repository, since  
520 the pre-processing can be conducted once and for  
521 all. However, it may be unsatisfactory in scenarios  
522 where users expect real-time uploads of new PDFs.  
523 2) The iterative chain-of-retrieval procedure incurs  
524 extra time delay compared to classic RAG. This is  
525 anticipated and Snell et al. (2024) points out that  
526 scaling inference computation is more effective  
527 than scaling model parameters. 3) Different from  
528 GraphRAG (Edge et al., 2024) which can pre-parse  
529 any free-form texts, NeuSym-RAG requires that  
530 the external documents exhibit implicit structures  
531 or layouts. And PDF files in AI research naturally  
532 conform to this constraint. Thus, integrating the  
533 two retrieval paradigms can maximize the utiliza-  
534 tion of the inherent advantages of semi-structured  
535 documents. 4) Currently, AIRQA-REAL are lim-  
536 ited to the domain of AI research papers. We are  
537 actively working on expanding our efforts to other  
538 vertical domains, including finance and law.

## 539 Ethics Statement

540 All papers utilized in the construction of our cu-  
541 rated Q&A dataset, AIRQA-REAL, are sourced  
542 from the websites ACL Anthology <sup>3</sup>, arXiv.org <sup>4</sup>,  
543 and OpenReview.net <sup>5</sup>, which are all licensed under  
544 the Creative Commons Attribution 4.0 Interna-  
545 tional License (CC BY 4.0). The use of these  
546 materials complies fully with their licensing terms,  
547 strictly for academic and research purposes. No  
548 private, sensitive, or personally identifiable infor-  
549 mation is used in this work. The study adheres to  
550 the ethical guidelines outlined by ACL, ensuring  
551 integrity, transparency, and reproducibility.

## 552 References

553 Amar Abane, Anis Bekri, and Abdella Battou. 2024.  
554 Fastrag: Retrieval augmented generation for semi-

<sup>3</sup>ACL Anthology: <https://aclanthology.org/>

<sup>4</sup>arXiv.org API: [https://info.arxiv.org/help/api\\_index.html](https://info.arxiv.org/help/api_index.html)

<sup>5</sup>OpenReview.net APIv2: <https://docs.openreview.net/reference/api-v2>

555 structured data. *arXiv e-prints*, pages arXiv–2411.

Anirudh Ajith, Mengzhou Xia, Alexis Chevalier, Tanya Goyal, Danqi Chen, and Tianyu Gao. 2024. Lit-search: A retrieval benchmark for scientific literature search.

Inc. Artifex Software. 2023. Pymupdf - a python binding for mupdf. <https://pymupdf.readthedocs.io/en/latest/>. Version 1.24.9, accessed on January 25, 2025.

Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. Semantic parsing on freebase from question-answer pairs. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1533–1544.

Asim Biswal, Liana Patel, Siddarth Jha, Amog Kamsetty, Shu Liu, Joseph E Gonzalez, Carlos Guestrin, and Matei Zaharia. 2024. Text2sql is not enough: Unifying ai and databases with tag. *arXiv preprint arXiv:2408.14717*.

Zheng Cai, Maosong Cao, Haojong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhao Fei, Yang Gao, Jiaye Ge, Chenya Gu, Yuzhe Gu, Tao Gui, Aijia Guo, Qipeng Guo, Conghui He, Yingfan Hu, Ting Huang, Tao Jiang, Penglong Jiao, Zhenjiang Jin, Zhikai Lei, Jiaxing Li, Jingwen Li, Linyang Li, Shuaibin Li, Wei Li, Yining Li, Hongwei Liu, Jiangning Liu, Jiawei Hong, Kaiwen Liu, Kuikun Liu, Xiaoran Liu, Chengqi Lv, Haijun Lv, Kai Lv, Li Ma, Runyuan Ma, Zerun Ma, Wenchang Ning, Linke Ouyang, Jiantao Qiu, Yuan Qu, Fukai Shang, Yunfan Shao, Demin Song, Zifan Song, Zhihao Sui, Peng Sun, Yu Sun, Huanze Tang, Bin Wang, Guoteng Wang, Jiaqi Wang, Jiayu Wang, Rui Wang, Yudong Wang, Ziyi Wang, Xingjian Wei, Qizhen Weng, Fan Wu, Yingtong Xiong, Chao Xu, Ruiliang Xu, Hang Yan, Yirong Yan, Xiaogui Yang, Haochen Ye, Huaiyuan Ying, Jia Yu, Jing Yu, Yuhang Zang, Chuyu Zhang, Li Zhang, Pan Zhang, Peng Zhang, Ruijie Zhang, Shuo Zhang, Songyang Zhang, Wenjian Zhang, Wenwei Zhang, Xingcheng Zhang, Xinyue Zhang, Hui Zhao, Qian Zhao, Xiaomeng Zhao, Fengzhe Zhou, Zaida Zhou, Jingming Zhuo, Yicheng Zou, Xipeng Qiu, Yu Qiao, and Dahua Lin. 2024. Internlm2 technical report.

Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan R. Routledge, and William Yang Wang. 2021. Finqa: A dataset of numerical reasoning over financial data. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7–11 November, 2021*, pages 3697–3711. Association for Computational Linguistics.

Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A. Smith, and Matt Gardner. 2021. A dataset

612	of information-seeking questions and answers anchored in research papers.	In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021</i> , pages 4599–4610. Association for Computational Linguistics.	667
613		Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. Query rewriting in retrieval-augmented large language models.	668
614		In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 5303–5315.	669
615		670	
616			671
617			
618			
619	Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From local to global: A graph RAG approach to query-focused summarization.	<i>CoRR</i> , abs/2404.16130.	672
620			673
621		<i>DBI Package for the DuckDB Database Management System</i> . R package version 1.1.3.9017,	674
622		<a href="https://github.com/duckdb/duckdb-r">https://github.com/duckdb/duckdb-r</a> .	675
623			
624	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning.	<i>arXiv preprint arXiv:2501.12948</i> .	676
625			677
626			678
627			679
628			680
629	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training.	In <i>International conference on machine learning</i> , pages 3929–3938. PMLR.	681
630			682
631			
632			
633	Xueyu Hu, Ziyu Zhao, Shuang Wei, Ziwei Chai, Qianli Ma, Guoyin Wang, Xuwu Wang, Jing Su, Jingjing Xu, Ming Zhu, Yao Cheng, Jianbo Yuan, Jiwei Li, Kun Kuang, Yang Yang, Hongxia Yang, and Fei Wu. 2024. Infiagent-dabench: Evaluating agents on data analysis tasks.	In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	683
634			684
635			685
636			686
637			687
638			688
639			689
640			
641	Yizheng Huang and Jimmy Huang. 2024. A survey on retrieval-augmented text generation for large language models.	<i>CoRR</i> , abs/2404.10981.	690
642			691
643			692
644	Hugging Face Team. 2023. Llama 3.3-70b-instruct model.	<a href="https://huggingface.co/meta-llama/Llama-3.3-70B-Instruct">https://huggingface.co/meta-llama/Llama-3.3-70B-Instruct</a> .	693
645	Accessed: 2023-02-11.		
646			
647			
648	Yoonjoo Lee, Kyungjae Lee, Sungyun Park, Dasol Hwang, Jaehyeon Kim, Hong-in Lee, and Moontae Lee. 2023. Qasa: advanced question answering on scientific articles.	In <i>Proceedings of the 40th International Conference on Machine Learning, ICML’23. JMLR.org</i> .	694
649			695
650			696
651			697
652			698
653			699
654	Chuhan Li, Ziyao Shangguan, Yilun Zhao, Deyuan Li, Yixin Liu, and Arman Cohan. 2024a. M3SciQA: A multi-modal multi-document scientific QA benchmark for evaluating foundation models.	In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 15419–15446, Miami, Florida, USA. Association for Computational Linguistics.	700
655			701
656			702
657			703
658			
659			
660			
661	Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhu Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. 2024b. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls.	<i>Advances in Neural Information Processing Systems</i> , 36.	704
662			705
663			706
664			
665			
666			
667	Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy.	In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9248–9274, Singapore. Association for Computational Linguistics.	715
668			716
669			717
670			718
671			719
672			720
673			721

- 722 Shruti Singh, Nandan Sarkar, and Arman Cohan. 2024. 776  
 723 [SciDQA: A deep reading comprehension dataset over](#) 777  
 724 [scientific papers.](#) In *Proceedings of the 2024 Conference* 778  
 725 [on Empirical Methods in Natural Language Processing.](#) 779  
 726
- 727 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. Scaling llm test-time compute optimally 780  
 728 can be more effective than scaling model parameters. 781  
 729 *arXiv preprint arXiv:2408.03314.*
- 730
- 731 Qwen Team. 2025. [Qwen2.5-vl.](#)
- 732 Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, 782  
 733 and Ashish Sabharwal. 2023. Interleaving retrieval 783  
 734 with chain-of-thought reasoning for knowledge- 784  
 735 intensive multi-step questions. In *Proceedings of the* 785  
 736 *61st Annual Meeting of the Association for Compu-*  
 737 *tational Linguistics (Volume 1: Long Papers)*, pages  
 738 10014–10037.
- 739 Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, 782  
 740 Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan 783  
 741 Qu, Fukai Shang, Bo Zhang, Liqun Wei, Zhihao Sui, 784  
 742 Wei Li, Botian Shi, Yu Qiao, Dahua Lin, and Conghui 785  
 743 He. 2024a. [Mineru: An open-source solution for](#)  
 744 [precise document content extraction.](#)
- 745 Jiaoguo Wang, Xiaomeng Yi, Rentong Guo, Hai Jin, 782  
 746 Peng Xu, Shengjun Li, Xiangyu Wang, Xiangzhou 783  
 747 Guo, Chengming Li, Xiaohai Xu, et al. 2021. Milvus: 784  
 748 A purpose-built vector data management system. In 785  
 749 *Proceedings of the 2021 International Conference on*  
 750 *Management of Data*, pages 2614–2627.
- 751 Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan 782  
 752 Yang, and Ming Zhou. 2020. [Minilm: Deep self-](#)  
 753 [attention distillation for task-agnostic compression](#)  
 754 [of pre-trained transformers.](#) In *Advances in Neural*  
 755 *Information Processing Systems 33: Annual Confer-*  
 756 *ence on Neural Information Processing Systems 2020,*  
 757 *NeurIPS 2020, December 6-12, 2020, virtual.*
- 758 Zhongyuan Wang, Richong Zhang, Zhijie Nie, and Jaein 782  
 759 Kim. 2024b. [Tool-assisted agent on SQL inspec-](#)  
 760 [tion and refinement in real-world scenarios.](#) *CoRR*,  
 761 [abs/2408.16991.](#)
- 762 Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas 782  
 763 Muennighoff. 2023. [C-pack: Packaged resources](#)  
 764 [to advance general chinese embedding.](#)
- 765 Fangyuan Xu, Weijia Shi, and Eunsol Choi. 2023. Re- 782  
 766 comp: Improving retrieval-augmented lms with  
 767 compression and selective augmentation. *arXiv preprint*  
 768 *arXiv:2310.04408.*
- 769 Hongshen Xu, Zichen Zhu, Situo Zhang, Da Ma, Shuai 782  
 770 Fan, Lu Chen, and Kai Yu. 2024. Rejection im-  
 771 proves reliability: Training llms to refuse unknown  
 772 questions using rl from knowledge feedback. *arXiv*  
 773 *preprint arXiv:2403.18349.*
- 774 S. Yang. 2023. [Advanced rag 01: Small-to-big retrieval.](#)  
 775 Accessed on 2024-12-11.

786           **A AIRQA-REAL Dataset**

787           **A.1 Data Format**

788       In this section, we briefly introduce the data format of our AIRQA-REAL dataset. Each data instance is  
789       represented as a JSON dictionary which contains the following fields:

- 790       • **uuid**: Globally unique uuid of the current task example.
- 791       • **question**: The user question about the given papers.
- 792       • **answer\_format**: The requirements for the output format of LLMs, e.g., a Python list of text strings or  
793        a single float number, such that we can evaluate the answer conveniently.
- 794       • **tags**: A list of tags denoting different types, categories or genres. Feasible tags include [“single”,  
795        “multiple”, “retrieval”, “text”, “image”, “table”, “formula”, “metadata”, “subjective”, “objective”].
- 796       • **anchor\_pdf**: A list of PDF uuids that are directly related to or explicitly mentioned in the question.
- 797       • **reference\_pdf**: A list of PDF uuids that may or may not help answer the question.
- 798       • **conference**: A list of conference names plus year. This field is only useful in the “*paper retrieval*”  
799        setting to limit the scope of search space, as shown in the task example of Listing 4.
- 800       • **evaluator**: A dictionary containing 2 fields, `eval_func` and `eval_kwargs`, which defines how to  
801        evaluate the model outputs. Concretely,
  - 802           ■ the “`eval_func`” field defines the name of our customized Python function (or metric) which is  
803            used to compare the predicted result and the expected ground truth;
  - 804           ■ the “`eval_kwargs`” field defines the arguments for the corresponding evaluation function, which  
805            usually contain the gold or reference answer and other optional parameters.

806       For example, in Listing 3, we utilize the function “`eval_structured_object_exact_match`” to  
807       evaluate the LLMs output according to the given gold answer [“SCG-NLI”, “false”].

```
808       1 {  
809       2     "uuid": "927ff9af-42f7-5216-a6f9-f106e8ff6759",  
810       3     "question": "On the HEML sentence level with AUC metric, which baseline  
811            ↳ outperforms MIND on specific conditons? Is it the best variant according  
812            ↳ to the paper that proposed that baseline?",  
813       4     "answer_format": "Your answer should be a Python list of two strings. The first  
814            ↳ string is the name of the baseline (with variant) that outperforms MIND,  
815            ↳ as proposed in the anchor PDF. The second string is either `true` or  
816            ↳ `false`.",  
817       5     "tags": [  
818       6        "multiple",  
819       7        "text",  
820       8        "table",  
821       9        "objective"  
822       10      ],  
823       11      "anchor_pdf": [  
824       12        "621d42a1-dbab-5003-b7c5-625335653001"  
825       13      ],  
826       14      "reference_pdf": [  
827       15        "ab661558-432d-5e5e-b49c-a3660a40986e",  
828       16        "1a21b653-3db0-55e8-9d34-8b6cd3dcbea",  
829       17        "85111b8b-4df0-5a9a-8d11-a7ae12eebcf6",  
830       18        "0597ce2b-cd8c-5b5b-b692-e8042d8548de",  
831       19        "6df0f3f3-e2e1-5d7a-9d70-3114ceac5939",  
832       20        "02f7ffff5-cec7-5ac8-a037-f5eb117b9547"  
833       21      ],  
834       22      "conference": [],  
835       23      "evaluator": {  
836       24        "eval_func": "eval_structured_object_exact_match",
```

```

25     "eval_kwargs": {
26         "gold": [
27             "SCG-NLI",
28             "false"
29         ],
30         "ignore_order": false,
31         "lowercase": true
32     }
33 }
34 }
```

838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848

Listing 3: A multi-doc task example of JSON format from the dataset AIRQA-REAL.

```

1 {
2     "uuid": "4b4877cd-4cdc-5d52-ac20-edfaa6dd7e32",
3     "question": "Is there any paper leverages knowledge distillation of language
4         ↪ models for textual out-of-distribution detection or anomaly detection?",
5     "answer_format": "Your answer should be the title of the paper WITHOUT ANY
6         ↪ EXPLANATION.",
7     "tags": [
8         "retrieval",
9         "text",
10        "objective"
11    ],
12    "anchor_pdf": [],
13    "reference_pdf": [],
14    "conference": [
15        "acl2023"
16    ],
17    "evaluator": {
18        "eval_func": "eval_paper_relevance_with_reference_answer",
19        "eval_kwargs": {
20            "reference_answer": "Multi-Level Knowledge Distillation for
21                ↪ Out-of-Distribution Detection in Text"
22        }
23    }
24 }
```

849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873

Listing 4: A paper retrieval task example of JSON format from the dataset AIRQA-REAL.

## A.2 Data Examples

Table 7 shows 3 examples in AIRQA-REAL with different task types, namely single-doc details, multi-doc analysis, and paper retrieval. Different numbers of PDF documents are involved in different task types. For instance, single-doc examples only specify one paper in its “anchor\_pdf” field, while multi-doc examples must mention multiple papers in the “anchor\_pdf” or “reference\_pdf” fields. And for paper retrieval, we relax the search scope to the entire conference venue, such as “acl2023” in Listing 4.

Category	Question	Answer Format
single	On the ALFWorld dataset experiments, how much did the success rate improve when the authors used their method compared to the original baseline model?	Your answer should be a floating-point number with one decimal place.
multiple	I would like to reproduce the experiments of KnowGPT, could you please provide me with the websites of the datasets applied in the experiment?	Your answer should be a Python list of 3 strings, the websites. Note that you should provide the original URL as given in the papers that proposed the datasets.
retrieval	Find the NLP paper that focuses on dialogue generation and introduces advancements in the augmentation of one-to-many or one-to-one dialogue data by conducting augmentation within the semantic space.	Your answer should be the title of the paper WITHOUT ANY EXPLANATION.

Table 7: Demonstration examples of different task types from the dataset AIRQA-REAL, where `single`, `multiple` and `retrieval` represent single-doc details, multi-doc analysis and paper retrieval, respectively.

875  
876  
877  
878  
879

### A.3 Dataset Statistics

In total, the entire AIRQA-REAL contains 6797 PDF documents and 553 questions with respect to AI research papers mostly in year 2023 and 2024. Among them, 6384 PDFs are published papers from ACL 2023 and ICLR 2024, while the PDFs left are also from publicly available websites arXiv.org and OpenReview.net. Table 8, Figure 7 and Figure 8 show the statistics of the dataset AIRQA-REAL. Note that, one task example can be labeled with more than one category tag (i.e., text, table, image, formula, and metadata). As for the 553 examples, 240 are converted from the LitSearch (Ajith et al., 2024) benchmark (a paper retrieval task set). And only examples in LitSearch whose ground truth belongs to ACL 2023 or ICLR 2024 venues are selected and revised to be compatible with our AIRQA-REAL data format. Correspondingly, we set the value of field “conference” (defined in App. A.1) to be “acl2023” or “iclr2024” in order to restrict the paper search space.

Data Splits	Number	Ratio(%)
single	159	29
multiple	150	27
retrieval	244	44
text	470	85
table	101	18
image	42	8
formula	36	7
metadata	22	4
objective	446	81
subjective	107	19
Overall	553	100

Table 8: The number and ratio of different data splits on the dataset AIRQA-REAL.

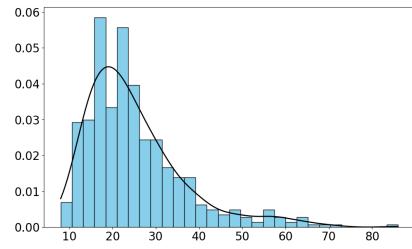


Figure 7: Frequency distribution of question lengths.

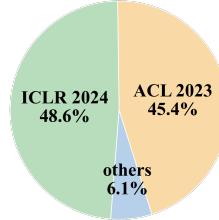


Figure 8: Conference distribution of the papers used.

### A.4 Other Full-PDF-based Academic Research Q&A Benchmarks

Apart from AIRQA-REAL, we also convert two other existing Q&A benchmarks M3SciQA (Li et al., 2024a) and SciDQA (Singh et al., 2024) in academic research to validate the universality of our NeuSym-RAG framework, which gives us 492 and 2937 more test samples, respectively. Specifically,

**M3SciQA** The answers of the test set on M3SciQA are not publicly available yet till 2025-04-11. Thus, we only utilize the complete 452 labeled samples in the validation set as the test suite which are licensed under Creative Commons Attribution 4.0. Notice that, each question in M3SciQA involves multiple papers. To obtain the final answer, agents need to firstly find the correct paper which belongs to the references of the anchor paper. Since our constructed database contains the entire paper set on M3SciQA, we only provide the title of the anchor PDF without referenced paper titles in the prompt to make the testing scenario more realistic. And agents need to find the target paper title by itself in the reference section. PDFs of all relevant papers can be downloaded from the arXiv.org website. As for evaluation, we follow the official evaluation metric <sup>6</sup> and leverage large language model based assessment with exactly the same prompt, model (gpt-4-0125-preview) and temperature (0.0) to determine whether the final long-form answer is correct.

**SciDQA** We convert the original 2937 test samples which are licensed under Open Data Commons Attribution License (ODC-By) v1.0. The relevant paper PDFs can all be downloaded from OpenReview.net and we categorize each input question as “table”, “image” and “formula” based on heuristic rules (e.g., whether a specific keyword such as “equation” is mentioned in the question or ground

<sup>6</sup>Evaluation for M3SciQA: [https://github.com/yale-nlp/M3SciQA/blob/main/src/evaluate\\_detail.py](https://github.com/yale-nlp/M3SciQA/blob/main/src/evaluate_detail.py)

truth). As for evaluation, we adopt the official LLM-based judgement and answer extraction scripts <sup>7</sup> except that the large language model is replaced with gpt-4o-mini from Llama-3.1-70B-Instruct due to its performance being more closely aligned with humans according to Figure 6.

Two separate evaluation functions are encapsulated for these two benchmarks (see Table 12). The authors declare that the usage of these two existing benchmarks strictly obeys the aforementioned licenses.

## B Supplementary Experiments and Settings

For NeuSym-RAG, we also investigate different formats of the action and observation spaces (exemplified in App. F.2.2 and App. F.2.3). According to Table 9: 1) The action format has a more substantial impact on the results compared to the observation format. 2) Specifically, both LLMs perform the best when the action format is “*markdown*”, which resembles the function calling fashion in Python. Another possible explanation is that the other 3 action formats generally impose stricter formatting requirements. 3) Different LLMs may prefer different observation formats on account of their training corpus. But the gaps, especially for the top-performing 2 choices, are relatively small.

Type	Format	Overall(%)	Type	Format	Overall(%)
GPT-4o-mini-2024-07-18			Qwen2.5-72B-Instruct		
action	<b>markdown</b>	<b>40</b>	action	<b>markdown</b>	<b>28</b>
	json	35		json	27
	xml	28		xml	19
	yaml	32		yaml	26
observation	<b>json</b>	<b>40</b>	observation	json	28
	markdown	31		markdown	27
	html	31		html	26
	string	39		<b>string</b>	<b>30</b>

Table 9: Ablation study on different action and observation formats on a subset of AIRQA-REAL dataset.

In Table 10, we present some statistics of iterative methods on different models in our experiments. In general, GPT-series LLMs complete the task with fewer interaction rounds, indicating that closed-source LLMs exhibit stronger capabilities and confidence. Using the Qwen2.5-72B-Instruct LLM, the NeuSym-RAG approach completes the task with more interaction rounds than the other two iterative methods, since the NeuSym-RAG approach can integrate the context from both retrieval paradigms.

Model	RAG Method	# Interaction(s)	# Prompt Token(s)	# Completion Token(s)	Time (s)	Cost (\$)
Qwen2.5-72B-Instruct	NeuSym-RAG	4.45	36521	832	40	-
	Iterative Sym-RAG	5.16	43548	2730	105	-
	Iterative Neu-RAG	4.19	48600	2395	94	-
	NeuSym-RAG	5.26	58262	1987	83	-
GPT-4o-mini	NeuSym-RAG	3.59	29306	518	20	0.0059
GPT-4V	NeuSym-RAG	3.61	54657	1098	17	0.1464

Table 10: Statistics of the number of interaction(s), accumulated prompt / completion token(s), time consumption, and LLM cost per sample with different models and RAG methods on 100 samples from AIRQA-REAL.

Table 11 demonstrates the performance of NeuSym-RAG based on the Qwen2.5-72B-Instruct LLM with different text modality collections (*i.e.*, encoding models) in the vectorstore. The experimental results uncover that a single BM25 text collection suffices to attain the best performance. We speculate that there might be two possible reasons for this observation: 1) existing LLMs are still not capable of exploring retrieved context from multiple encoding models and then selecting the optimal one. This behavioral pattern requires profound reasoning abilities. And we attempt to verify this hypothesis by adopting supervised fine-tuning with agent trajectories in our future work. 2) Another potential insight is that, for different encoding models, the best practice is to directly choose the highest-performing one

<sup>7</sup>Evaluation for SciDQA: <https://github.com/yale-nlp/SciDQA/tree/main/src/evals>

instead of leaving the burden to the agent. And in our preliminary experiments, we also find similar phenomena that instead of creating different DB tables to store the parsed document content with different PDF parsing tools, directly choosing the top-performing tools, that is PyMuPDF ([Artifex Software, 2023](#)) and MinerU ([Wang et al., 2024a](#)) in § 2.2, and using a single DB table achieves the best performance.

BM25	all-MiniLM-L6-v2	bge-large-en-v1.5	text	table	image	formula	metadata	Overall(%)
✓	✗	✗	48.84	35.00	38.89	61.90	45.45	46
✗	✓	✗	41.86	25.00	27.78	33.33	27.27	33
✗	✗	✓	39.53	25.00	16.67	52.38	54.55	37
✓	✓	✗	44.19	30.00	27.78	52.38	27.27	40
✓	✗	✓	41.86	20.00	33.33	42.86	27.27	35
✗	✓	✓	34.88	10.00	16.67	47.62	45.45	30
✓	✓	✓	39.53	20.00	16.67	47.62	45.45	35

Table 11: Performance of NeuSym-RAG on AIRQA-REAL dataset with different text encoding model(s) applied. The experiments are conducted using open-source LLM Qwen2.5-72B-Instruct.

## C Agent Baselines

In this section, we elaborate the implementation of each RAG method utilized in Table 3.

- **Classic RAG** fetches question-related chunks from the vectorstore and directly provides them as the context for LLMs to answer the question. The chunk size is fixed to 512 tokens using the RecursiveCharacterTextSplitter from langchain<sup>8</sup> and the top K size is set to 4. In other words, we pre-fetch the content of column “chunks.text\_content” based on the raw input question and insert them as the context in the prompt. The default text embedding model is fixed to all-MiniLM-L6-v2.
- **Iterative Classic RAG** enables LLMs to repeatedly alter their query texts and iteratively retrieve chunks until the answer can be obtained. But the view is always fixed to column “chunks.text\_content”.
- **Two-stage Neu-RAG** splits the task into two stages. At stage one, LLMs predict a RETRIEVEFROMVECTORSTORE action with only one chance. But they can predict the parameters in that action, e.g., the query text, the column to choose, the embedding collection, and the top K size. While at stage two, agents must output the final answer with retrieved context.
- **Iterative Neu-RAG** allows for trial and error. LLMs can predict multiple parameterized actions to retrieve from the vectorstore until the interaction trajectory suffices to answer the question. And we incorporate another two useful actions VIEWIMAGE and CALCULATEEXPR (formally defined in § 2.4.1 and App. F.2.1) into the multi-turn interaction.
- **Two-stage Sym-RAG** requires LLMs to generate a SQL query first. After SQL execution upon the backend database, they must predict the answer using the retrieved context, with only one chance.
- **Iterative Sym-RAG** is the multi-round version, where LLMs can iteratively interact with the symbolic database and try different SQLs to better conclude the final answer. The additional two action types VIEWIMAGE and CALCULATEEXPR are also included in the action space like Iterative Neu-RAG.
- **Graph-RAG** ([Edge et al., 2024](#)) is implemented by following the official guideline of library graphrag<sup>9</sup>. And we adopt the local search mode since it performs better than global in our pilot study.
- **Hybrid-RAG** ([Sarmah et al., 2024](#)) can be considered as a two-stage structured baseline, where the LLM agent firstly determines which action (RETRIEVEFROMVECTORSTORE or RETRIEVEFROMDATABASE) to use. After predicting the parameters of the action, and fetching the context from either the vectorstore or database, it needs to generate the final answer with only one chance.

<sup>8</sup>Langchain text splitter: [https://python.langchain.com/docs/how\\_to/recursive\\_text\\_splitter/](https://python.langchain.com/docs/how_to/recursive_text_splitter/).

<sup>9</sup>GraphRAG website: [https://microsoft.github.io/graphrag/get\\_started/](https://microsoft.github.io/graphrag/get_started/).

## D Evaluation Metrics

The detailed information of all 18 evaluation functions we design is presented in Table 12.

969

## E Database Schema and Encodable Columns in Vectorstore

The complete database schema to store parsed PDF content is visualized in Figure 9, with free online tool drawSQL<sup>10</sup>. All encodable columns are listed in Table 13. These columns provide various perspectives towards interpreting the PDF content. And these different views are inherently connected by the sub-graph of the original database schema (as illustrated in Figure 10).

972

973

974

975

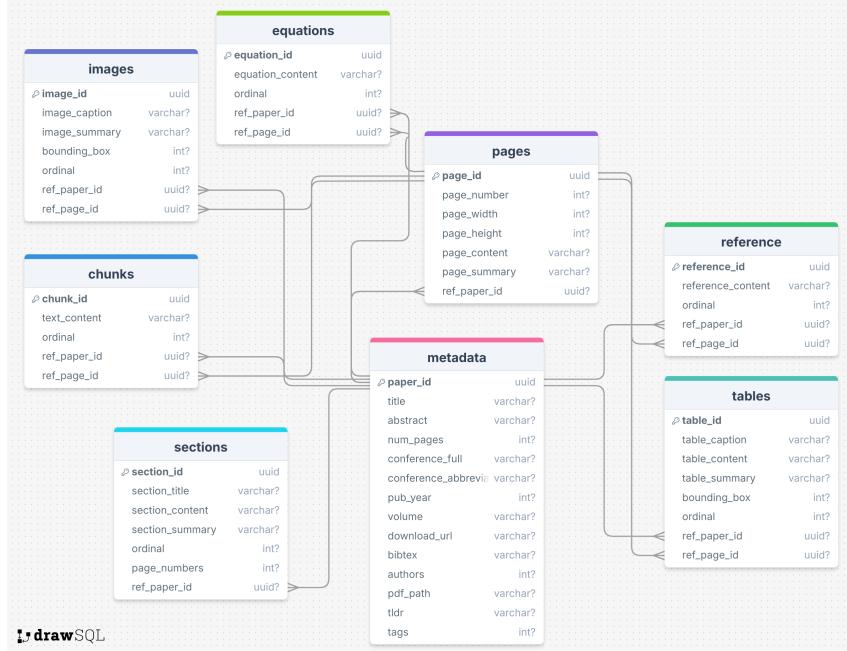


Figure 9: The complete and universal database schema to store the parsed elements of each PDF file. Note that, since the visualization tool drawSQL cannot display “ARRAY” data types, the actual data types of the columns “*bounding\_box*”, “*authors*”, “*tags*”, and “*page\_numbers*” are INT[4], VARCHAR[], VARCHAR[], and INT[], respectively.

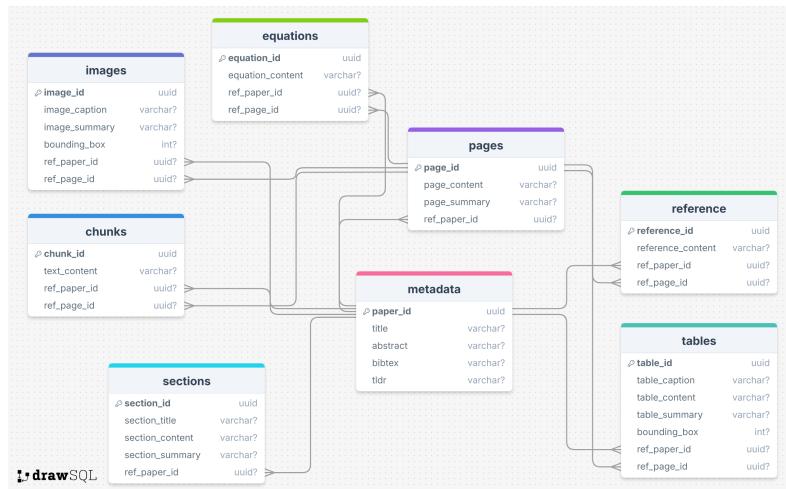


Figure 10: All encodable columns which are inherently connected by the schema sub-graph. Note that, since the visualization tool drawSQL cannot display “ARRAY” data types, the actual data types of columns “*images.bounding\_box*” and “*tables.bounding\_box*” are both INT[4].

<sup>10</sup><https://drawsql.app/>

<b>Eval Type</b>	<b>Category</b>	<b>Function</b>	<b>Description</b>
<i>objective</i>	match	eval_bool_exact_match	Evaluate the output against the answer using exact boolean match.
		eval_float_exact_match	Evaluate the output against the answer using exact float match with variable precision or tolerance.
		eval_int_exact_match	Evaluate the output against the answer using exact integer match.
		eval_string_exact_match	Evaluate the output against the answer using exact string match.
		eval_structured_object_exact_match	Evaluate the output against the answer recursively by parsing them both as Python-style lists or dictionaries.
	set	eval_element_included	Evaluate whether the output is included in the answer list.
		eval_element_list_included	Evaluate whether each element in the output list is included in the answer list.
		eval_element_list_overlap	Evaluate whether the output list overlaps with the answer list.
	retrieval	eval_paper_relevance_with_reference_answer	Evaluate whether the retrieved paper is the same as the reference answer.
<i>subjective</i>	semantic	eval_reference_answer_with_llm	Evaluate the output against the reference answer using LLMs.
		eval_scoring_points_with_llm	Evaluate whether the scoring points are all mentioned in the output using LLMs.
		eval_partial_scoring_points_with_llm	Evaluate whether the scoring points are partially mentioned in the output using LLMs.
	formula	eval_complex_math_formula_with_llm	Evaluate the mathematical equivalence between the output and the answer formatted in Latex using LLMs.
	logical	eval_conjunction	Evaluate the conjunction of multiple evaluation functions. The output passes the evaluation if and only if all the elements in the output pass the corresponding sub-evaluations.
		eval_disjunction	Evaluate the disjunction of multiple evaluation functions. The output passes the evaluation if and only if at least one of the element in the output passes the corresponding sub-evaluation.
		eval_negation	Evaluate the negation of an evaluation function. The output passes the evaluation if and only if it doesn't pass the original evaluation function.
	others	eval_m3sciqqa	Evaluate examples in dataset M3SciQA with the encapsulated original LLM-based function.
		eval_scidqa	Evaluate examples in dataset SciDQA with the encapsulated original LLM-based function.

Table 12: The checklist of the 18 used evaluation functions, including their categories, names, and descriptions.

<b>Table</b>	<b>Column</b>	<b>Description</b>
metadata	title	The title of this paper.
	abstract	The abstract of this paper.
	bibtex	The bibtex of this paper.
	tldr	A brief summary of the paper's main idea or findings generated by LLM based on title and abstract.
pages	page_content	The content of the page.
	page_summary	A brief summary of the page content, generated by LLM, focusing on key information and describing the page content.
images	image_caption	The caption of this image, empty string if not found.
	image_summary	A brief summary of the image, generated by LLM, focusing on key information and describing the image.
	bounding_box	The bounding box of the figure in the format $[x_0, y_0, w, h]$ , where $(x_0, y_0)$ represents the coordinates of the top-left corner and $(w, h)$ represents the width and height which are used to determine the shape of the rectangle. The cropped image is encoded.
tables	table_caption	Caption of the table, showing key information of the table.
	table_content	The content of the table in html format.
	table_summary	A brief summary of the table content generated by LLM, focusing on key information and describing the table content.
	bounding_box	The bounding box of the table in the format $[x_0, y_0, w, h]$ , where $(x_0, y_0)$ represents the coordinates of the top-left corner and $(w, h)$ represents the width and height. The cropped image is encoded.
sections	section_title	The title of the current section.
	section_content	The text content of the current section.
	section_summary	A brief summary of the section content generated by LLM, focusing on key information and describing the section content.
chunks	text_content	The text content of the current chunk.
equations	equation_content	Content of the equation in latex format.
reference	reference_content	Text content of each reference.

Table 13: The checklist of all encodable columns with their descriptions.

976

## F Prompt Template

This section presents the detailed structure of the prompts used in various agent baselines (Figure 4), outlining the components that shape the agent's interactions and reasoning process. The overall prompt template includes the following five key components.

### Overall Prompt Composition

**[System Prompt]:** Defines the agent's role and describes the task to tackle, clarifying the context of the task.

**[Action and Observation Space Prompt]:** Defines the list of all feasible actions that the agent can take, including the action description, observation space, syntax and parameters, and use cases for each type. The detailed specification varies depending on the action format.

**[Interaction Framework Prompt]:** Outlines the main interaction procedure and template.

**[Hint Prompt]:** Provides hints or suggestions to help the agent refine its planning and reasoning process.

**[Task Prompt]:** Defines the task input, usually including the input question, required answer format, retrieved context, database schema, vectorstore schema and operators.

980

### E.1 System Prompt

The system prompt is used to define the ultimate goal that agents should achieve, the environment (if exists) that agents can interact with and the expected agent behavior to conduct. Notice that, the prompts for different agent baselines differ from each other.

#### System Prompts for Different Agent Baselines

##### For NeuSym-RAG:

You are an intelligent agent with expertise in **retrieving useful context from both the DuckDB database and the Milvus vectorstore through SQL execution and similarity search and answering user questions**. You will be given a natural language question concerning PDF files, along with the schema of both the database and the vectorstore. Your ultimate goal is to answer the input question with pre-defined answer format. The DuckDB database contains all parsed content of raw PDF files, while the Milvus vectorstore encodes specific column cells from the database as vectors. You can predict executable actions, interact with the hybrid environment (including database and vectorstore) across multiple turns, and retrieve necessary context until you are confident in resolving the question.

##### ## Task Description

Each input task consists of the following parts:

[Question]: A natural language question from the user regarding PDF files, e.g., Is there any ...?

[Answer Format]: Specifies the required format of the final answer, e.g., the answer is Yes or No without punctuation.

[Database Schema]: A detailed serialized schema of the DuckDB database for reference when generating SQL queries. It includes 1) tables, 2) columns and their data types, 3) descriptions for these schema items, and 4) primary key and foreign key constraints.

[Vectorstore Schema]: A detailed serialized schema of the Milvus vectorstore for reference when generating executable retrieval actions with specific parameters. It includes 1) collections, 2) fields, 3) encodable (table, column) pairs in the relational database where the vectorized content originates, and 4) grammar for valid filter rules.

##### For Classic RAG:

You are intelligent agent who is expert in **answering user questions** based on the retrieved context. You will be given a natural language question concerning a PDF file, and your task is to answer the input question with predefined output format using the relevant information.

##### For Two-stage Neu-RAG:

**[Stage 1]** You are intelligent agent who is expert in **predicting a well-formed retrieval action** to search useful information to answer the user question. You will be given a natural language question concerning a PDF file and a vectorstore schema which defines all usable collections and fields in them. The vectorized contents in the vectorstore all come from cell values in another relational database which stores the parsed content of the PDF files. And your task is to predict a parametrized retrieval action to find useful information based on vector similarity search. Please refer to the concrete vectorstore schema to produce a valid retrieval action.

**[Stage 2]** You are intelligent agent who is expert in **answering user question** given the retrieved context. You will be given a natural language question concerning a PDF file and the retrieved context. Your task is to predict the final answer based on given question and context. Please refer to the answer format to produce the valid answer.

985

## System Prompt

### For Iterative Classic RAG and Iterative Neu-RAG:

You are intelligent agent who is expert in **retrieving useful context from the vectorstore based on similarity search** and **answering user questions**. You will be given a natural language question concerning a PDF file and a vectorstore schema of Milvus, and your ultimate task is to answer the input question with pre-defined output format. The Milvus vectorstore encodes various context from the parsed PDF in multi-views. You can predict executable actions, interact with the vectorstore in multiple turns, and retrieve desired context to help you better resolve the question.

#### ## Task Description

Each input task consists of the following parts:

[Question]: A natural language question from the user regarding PDF files, e.g., Is there any ...?

[Answer Format]: Specifies the required format of the final answer, e.g., the answer is “Yes” or “No” without punctuation.

[Vectorstore Schema]: A detailed serialized schema of the Milvus vectorstore for reference when generating executable retrieval actions with specific parameters. It includes 1) collections, 2) fields, 3) encodable (table, column) pairs in the relational database where the vectorized content originates, and 4) grammar for valid filter rules.

### For Two-stage Sym-RAG:

[Stage 1] You are intelligent agent who is expert in **writing SQL programs** to retrieve useful information. You will be given a natural language question concerning a PDF file and a database schema which stores the parsed PDF content, and your task is to predict SQL to retrieve content from the database. Please refer to the concrete database schema to produce the valid SQL.

[Stage 2] the same as method Two-stage Neu-RAG

### For Iterative Sym-RAG:

You are intelligent agent who is expert in **leveraging SQL programs to retrieve useful information** and **answer user questions**. You will be given a natural language question concerning a PDF file and a database schema of DuckDB which stores the parsed PDF content, and your ultimate task is to answer the input question with predefined output format. You can predict intermediate SQLs, interact with the database in multiple turns, and retrieve desired information to help you better resolve the question.

#### ## Task Description

Each input task consists of the following parts:

[Question]: A natural language question from the user regarding PDF files, e.g., Is there any ...?

[Answer Format]: Specifies the required format of the final answer, e.g., the answer is Yes or No without punctuation.

[Database Schema]: A detailed serialized schema of the DuckDB database for reference when generating SQL queries. It includes 1) tables, 2) columns and their data types, 3) descriptions for these schema items, and 4) primary key and foreign key constraints.

### For Hybrid-RAG:

[Stage 1] You are intelligent agent who is expert in predicting a well-formed retrieval action to search useful information to answer the user question. You will be given a natural language question concerning a PDF file, a database schema which stores the parsed PDF content, and a vectorstore schema which defines all usable collections and fields in them. The vectorized contents in the vectorstore all come from cell values in the database. And your task is to predict a parametrized retrieval action to find useful information. Please refer to the concrete schema to produce a valid retrieval action.

[Stage 2] the same as methods Two-stage Neu-RAG and Two-stage Sym-RAG

986

## F.2 Action and Observation Space Prompt

In total, there are 5 actions (see § 2.4.1) for the proposed NeuSym-RAG framework. We choose the RETRIEVEFROMVECTORSTORE action as an example, and serialize it in JSON format below:

987

988

989

### Action and Observation Space Prompt for NeuSym-RAG (JSON format)

#### ## Action and Observation Space

All allowable action types include [“RetrieveFromVectorstore”, “RetrieveFromDatabase”, “CalculateExpr”, “ViewImage”, “GenerateAnswer”]. Here is the detailed specification in JSON format for them:

#### ### Action Type

RetrieveFromVectorstore

#### ### Description

Given a query text, retrieve relevant context from the Milvus vectorstore. Please refer to the schema of different collections and fields for each stored data entry.

#### ### Observation

The observation space is the retrieved top-ranked entries from the Milvus vectorstore based on input parameters.

990

## Action and Observation Space Prompt for NeuSym-RAG (JSON format) – continued

### ### Syntax and Parameters (JSON Format)

```
{  
    "action_type": "RetrieveFromVectorstore",  
    "parameters": {  
        "query": {  
            "type": "str",  
            "required": true,  
            "description": "This query will be encoded and used to search for relevant context.  
                You can rephrase the user question to obtain a more clear and structured requirement."  
        },  
        "collection_name": {  
            "type": "str",  
            "required": true,  
            "description": "The name of the collection in the Milvus vectorstore to search for  
                relevant context. Please ensure the collection does exist in the vectorstore."  
        },  
        "table_name": {  
            "type": "str",  
            "required": true,  
            "description": "The table name is used to narrow down the search space. And it will  
                be added to the filter condition. Please ensure this table has encodable columns."  
        },  
        "column_name": {  
            "type": "str",  
            "required": true,  
            "description": "The column name is used to narrow down the search space. And it will  
                be added to the filter condition. Please ensure it is encodable in `table_name`."  
        },  
        "filter": {  
            "type": "str",  
            "required": false,  
            "default": "",  
            "description": "The filter condition to narrow down the search space. Please refer to  
                the syntax of filter rules. By default, it is empty. It is suggested to restrict  
                `primary_key`, `pdf_id`, or `page_number` to refine search results."  
        },  
        "limit": {  
            "type": "int",  
            "required": false,  
            "default": 5,  
            "description": "The number of top-ranked context to retrieve. Please ensure that it is a  
                positive integer. And extremely large limit values may be truncated."  
        }  
    }  
}
```

### ### Use Cases (JSON Format)

#### #### Case 1

Search the Mivlus collection `text_bm25_en`, which use BM25 sparse embeddings, with the filter condition `"table_name == 'chunks'` and `column_name == 'text_content'` and `pdf_id == '12345678'` and `page_number == 1`" to restrict the content source and return the top 10 relevant entries.

[Action]:

```
{"action_type": "RetrieveFromVectorstore", "parameters": {"query": "Does this paper discuss LLM-based  
agent on its first page?", "collection_name": "text_bm25_en", "table_name": "chunks",  
"column_name": "text_content", "filter": "pdf_id == '12345678' and page_number == 1", "limit": 10}}
```

#### #### Case 2

... more cases in JSON format ...

- - -

... specification for other types of actions, the prompt can be easily inferred ...

## F.2.1 Syntax and Parameters for Other Action Types (JSON Format)

This subsection formally describes the syntax and parameters for the other 4 action types, including RETRIEVEFROMDATABASE, CALCULATEEXPR, VIEWIMAGE, GENERATEANSWER.

```

1  {
2      "action_type": "RetrieveFromDatabase",
3      "parameters": {
4          "sql": {
5              "type": "str",
6              "required": true,
7              "description": "The concrete DuckDB SQL query to execute and retrieve
8                  ↪ results."
9          }
10     },
11     {
12         "action_type": "CalculateExpr",
13         "parameters": {
14             "expr": {
15                 "type": "str",
16                 "required": true,
17                 "description": "The expression to calculate, e.g., '13 * 42'."}
18         }
19     },
20     {
21         "action_type": "ViewImage",
22         "description": "You can retrieve the visual information of the paper by taking
23             ↪ this action. Please provide the paper id, the page number, and the
24                 ↪ optional bounding box.",
25         "observation": "The observation space is the image that you want to view. We
26             ↪ will show you the image according to your parameters. The error message
27                 ↪ will be shown if there is any problem with the image retrieval.",
28         "parameters": {
29             "paper_id": {
30                 "type": "str",
31                 "required": true,
32                 "description": "The paper id to retrieve the image."}
33         },
34         "page_number": {
35             "type": "int",
36             "required": true,
37             "description": "The page number (starting from 1) of the paper to
38                 ↪ retrieve the image."}
39     },
40         "bounding_box": {
41             "type": "List[float]",
42             "required": false,
43             "default": [],
44             "description": "The bounding box of the image to retrieve. The format
45                 ↪ is [x_min, y_min, delta_x, delta_y]. The complete PDF page will
46                 ↪ be retrieved if not provided."}
47     },
48 },
49     {
50         "action_type": "GenerateAnswer",
51         "parameters": {
52             "answer": {
53                 "type": "Any",
54                 "required": true,
55                 "description": "The final answer to the user question. Please adhere to
56                     ↪ the answer format for the current question."}
57     }
58 }
```

Listing 5: Syntax and Parameters of JSON Format for Other Action Types in NeuSym-RAG

1059 Note that, except for the Classic RAG baseline, which has no action space, the feasible action types  
1060 differ among the other baselines described in Figure 4:

- **Iterative Classic RAG:** it only accepts the RETRIEVEFROMVECTORSTORE and GENERATEANSWER actions. For the former one, we further fix the perspective to be “chunks.text\_content” and the collection to text embedding model all-MiniLM-L6-v2;
- **Two-stage Neu-RAG:** it only accepts action RETRIEVEFROMVECTORSTORE at the first stage;
- **Iterative Neu-RAG:** all available action types during the iterative neural retrieval includes [RETRIEVEFROMVECTORSTORE, CALCULATEEXPR, VIEWIMAGE, GENERATEANSWER];
- **Two-stage Sym-RAG:** it only accepts action RETRIEVEFROMDATABASE at the first stage;
- **Iterative Sym-RAG:** all available action types during the iterative symbolic retrieval includes [RETRIEVEFROMDATABASE, CALCULATEEXPR, VIEWIMAGE, GENERATEANSWER];
- **Hybrid-RAG:** it accepts both actions RETRIEVEFROMVECTORSTORE and RETRIEVEFROMDATABASE at the first stage, and no action at the second stage.

## 1072 F.2.2 Other Action Formats Apart From JSON

1073 This subsection introduces other serialized action formats apart from JSON, namely MARKDOWN, XML  
1074 and YAML. Take action RETRIEVEFROMVECTORSTORE as an example (others can be easily inferred):

### MARKDOWN Format

#### ### Syntax and Parameters (MARKDOWN Format)

```
RetrieveFromVectorstore(query: str, collection_name: str, table_name: str, column_name: str, filter: str = '', limit: int = 5)
    - query: str, required. The query text will be encoded and used to search for relevant context.
        You can rephrase the user question to obtain a more clear and structured requirement.
    - collection_name: str, required. The name of the collection in the Milvus vectorstore to
        search for relevant context. Please ensure the collection does exist in the vectorstore.
    - table_name: str, required. The table name is used to narrow down the search space. And it
        will be added to the filter condition. Please ensure this table has encodable columns.
    - column_name: str, required. The column name is used to narrow down the search space. And it
        will be added to the filter condition. Please ensure it is encodable in `table_name`.
    - filter: str, optional, default to ''. The filter condition to narrow down the search space.
        Please refer to the syntax of filter rules. By default, it is empty. It is suggested to
        restrict `primary_key`, `pdf_id`, or `page_number` to refine search results.
    - limit: int, optional, default to 5. The number of top-ranked context to retrieve. Please
        ensure that it is a positive integer. And extremely large limit values may be truncated.
```

#### ### Use Cases (MARKDOWN Format)

##### #### Case 1

Search the Mivlus collection text\_bm25\_en, which use BM25 sparse embeddings, with the filter condition “table\_name == ‘chunks’ and column\_name == ‘text\_content’ and pdf\_id == ‘12345678’ and page\_number == 1” to restrict the content source and return the top 10 relevant entries.

[Action]:

```
RetrieveFromVectorstore(query="Does this paper discuss LLM-based agent on its first page?", collection_name='text_bm25_en', table_name='chunks', column_name='text_content', filter="pdf_id == '12345678' and page_number == 1", limit=10)
```

##### #### Case 2

Perform a vector-based similarity search on all cell values from the ‘abstract’ column in the ‘metadata’ table in the database, using the MiniLM-L6-v2 sentence transformer embeddings. By default, the top 5 most relevant entries will be returned.

[Action]:

```
RetrieveFromVectorstore(query="Is there any work about the topic structured RAG?", collection_name='text_sentence_transformers_all_minilm_l6_v2', table_name='metadata', column_name='abstract')
```

##### #### Case 3

.. more cases in MARKDOWN format ...

## XML Format

### ### Syntax and Parameters (XML Format)

```
<action>
  <action_type>RetrieveFromVectorstore</action_type>
  <parameters>
    <query>
      <type>str</type>
      <required>true</required>
      <description>The query text will be encoded and used to search for relevant context. You can rephrase the original user question to obtain a more clear and structured query requirement.</description>
    </query>
    <collection_name>
      <type>str</type>
      <required>true</required>
      <description>The collection name in the Milvus vectorstore to search for relevant context. Please ensure the collection does exist in the vectorstore.</description>
    </collection_name>
    <table_name>
      <type>str</type>
      <required>true</required>
      <description>The table name is used to narrow down the search space. It will be added to the filter condition. Please ensure this table has encodable columns.</description>
    </table_name>
    <column_name>
      <type>str</type>
      <required>true</required>
      <description>The column name is used to narrow down the search space. It will be added to the filter condition. Please ensure it is encodable in `table_name`.</description>
    </column_name>
    <filter>
      <type>str</type>
      <required>false</required>
      <default></default>
      <description>The filter condition to narrow down the search space. Please refer to the syntax of filter rules. By default, it is empty. It is suggested to restrict `pdf_id`, `page_number`, or `primary_key` to refine search results.</description>
    </filter>
    <limit>
      <type>int</type>
      <required>false</required>
      <default>5</default>
      <description>The number of top-ranked context to retrieve. Please ensure that it is a positive integer. And extremely large limit values may be truncated.</description>
    </limit>
  </parameters>
</action>
```

### ### Use Cases (XML Format)

#### #### Case 1

Search the Milvus collection `text_bm25_en`, which use BM25 sparse embeddings, with the filter condition “`table_name == 'chunks'` and `column_name == 'text_content'` and `pdf_id == '12345678'` and `page_number == 1`” to restrict the content source and return the top 10 relevant entries.

[Action]:

```
<action><action_type>RetrieveFromVectorstore</action_type><parameters><query>Does this paper discuss LLM-based agent on its first page?</query><collection_name>text_bm25_en</collection_name><table_name>chunks</table_name><column_name>text_content</column_name><filter>pdf_id == '12345678' and page_number == 1</filter><limit>10</limit></parameters></action>
```

#### #### Case 2

.. more cases in XML format ...

## YAML Format

### ### Syntax and Parameters (YAML Format)

```
action_type: RetrieveFromVectorstore
parameters:
  query:
    type: str
    required: true
    description: The query text will be encoded and used to search for relevant context. You can rephrase the user question to obtain a more clear and structured requirement.
  collection_name:
    type: str
    required: true
    description: The name of the collection in the Milvus vectorstore to search for relevant context. Please ensure the collection does exist in the vectorstore.
  table_name:
    type: str
    required: true
    description: The table name is used to narrow down the search space. It will be added to the filter condition. Please ensure this table has encodable columns.
  column_name:
    type: str
    required: true
    description: The column name is used to narrow down the search space. It will be added to the filter condition. Please ensure it is encodable in `table_name`.
  filter:
    type: str
    required: false
    default: ''
    description: The filter condition to narrow down the search space. Please refer to the syntax of filter rules. By default, it is empty. It is suggested to restrict `pdf_id`, `page_number` or `primary_key` to refine search results.
  limit:
    type: int
    required: false
    default: 5
    description: The number of top-ranked context to retrieve. Please set it to a positive integer to limit the number of returned results. Extremely large limit values may be truncated.
```

### ### Use Cases (YAML Format)

#### #### Case 1

Search the Mivlus collection `text_bm25_en`, which use BM25 sparse embeddings, with the filter condition “`table_name == 'chunks'` and `column_name == 'text_content'` and `pdf_id == '12345678'` and `page_number == 1`” to restrict the content source and return the top 10 relevant entries.

#### [Action]:

```
action_type: RetrieveFromVectorstore
parameters:
  query: Does this paper discuss LLM-based agent on its first page?
  collection_name: text_bm25_en
  table_name: chunks
  column_name: text_content
  filter: pdf_id == '12345678' and page_number == 1
  limit: 10
```

#### #### Case 2

... more cases in YAML format ...

1077

1078

## F.2.3 Observation Format

In this subsection, we discuss different observation formats to organize the retrieved entries. The execution result of action `RETRIEVEFROMDATABASE` naturally forms a *table*. While for action `RETRIEVEFROMVECTORSTORE`, we extract the `text` field for text type and the `bounding_box` field for image type respectively, and also concatenate the top-K entries as a *table*. We support the following 4 observation

formats to serialize the table-style observation: MARKDOWN, JSON, STRING, and HTML. 1083  
Take the RETRIEVEFROMDATABASE action as an example, the SQL query to execute is: 1084

```
select title, pub_year from metadata where conference_abbreviation = 'ACL' limit 3; 1085
```

Then, the returned observation texts in different formats are: 1086

#### Observation in MARKDOWN Format

title	pub_year
ContraCLM: Contrastive Learning For Causal Language Model	2023
Mitigating Label Biases for In-context Learning	2023
mCLIP: Multilingual CLIP via Cross-lingual Transfer	2023

In total, 3 rows are displayed in MARKDOWN format.

1087

#### Observation in JSON Format

```
{"title": "ContraCLM: Contrastive Learning For Causal Language Model", "pub_year": 2023}  
{"title": "Mitigating Label Biases for In-context Learning", "pub_year": 2023}  
{"title": "mCLIP: Multilingual CLIP via Cross-lingual Transfer", "pub_year": 2023}
```

In total, 3 rows are displayed in JSON format.

1088

#### Observation in STRING Format

	title	pub_year
ContraCLM: Contrastive Learning For Causal Language Model		2023
Mitigating Label Biases for In-context Learning		2023
mCLIP: Multilingual CLIP via Cross-lingual Transfer		2023

In total, 3 rows are displayed in STRING format.

1089

#### Observation in HTML Format

```
<table border="1" class="dataframe">  
  <thead>  
    <tr style="text-align: right;">  
      <th>title</th>  
      <th>pub_year</th>  
    </tr>  
  </thead>  
  <tbody>  
    <tr>  
      <td>ContraCLM: Contrastive Learning For Causal Language Model</td>  
      <td>2023</td>  
    </tr>  
    <tr>  
      <td>Mitigating Label Biases for In-context Learning</td>  
      <td>2023</td>  
    </tr>  
    <tr>  
      <td>mCLIP: Multilingual CLIP via Cross-lingual Transfer</td>  
      <td>2023</td>  
    </tr>  
  </tbody>  
</table>
```

In total, 3 rows are displayed in HTML format.

1090

### F.3 Interaction Framework Prompt

This section describes the interaction framework for iterative retrieval methods. For Classic and Two-stage RAG baselines, this part is omitted. We follow the popular ReAct framework (Yao et al., 2023) to encourage stepwise thought process before predicting actions.

1091

1092

1093

1094

## Interaction Framework Prompt

### ## Interaction Framework

The main interaction procedure proceeds like this:

- - -

[Thought]: reasoning process, why to take this action

[Action]: which action to take, please strictly conform to the action specification

[Observation]: execution results or error message after taking the action

... more interleaved triplets of ([Thought], [Action], [Observation]) ...

[Thought]: reasoning process to produce the final answer

[Action]: the terminal action `GenerateAnswer`, there is no further observation

- - -

In general, the main interaction loop consists of an interleaved of triplets ([Thought], [Action], [Observation]), except the last `GenerateAnswer` action which does not have "[Observation]:". You need to predict the "[Thought]: ..." followed by the "[Action]: ..." for each turn, and we will execute your action in the environment and provide the "[Observation]: ..." for the previous action.

1095

## F.4 Hint Prompt

This type of prompt provides a list of hints to guide the interaction. It highlights best practices for information retrieval, iterative refinement, and sequential decision-making throughout the task-solving process. It is highly customizable and can be easily extended to accommodate LLM prediction errors. Therefore, we only provide one demonstration example for the complete NeuSym-RAG framework.

### Hint Prompt for NeuSym-RAG

#### ## Suggestions or Hints for Agent Interaction

1. Explore multiple retrieval strategies. For example:

- Experiment with different (table, column) pairs to extract diverse types of information.
- Query various embedding models (collections) to find the most relevant context.

2. Combine both structured and unstructured data. Concretely:

- Use SQL queries to retrieve precise facts and structured data. Pay special attention to morphological variations in cell values.
- Perform similarity searches in the vectorstore to capture semantic relationships and hidden insights.

3. Iterate and refine:

- If SQL execution result is not satisfactory, try alternative SQL queries to explore the database content carefully.
- If the vector-based neural retrieval is insufficient, try alternative approaches or parameter settings.
- Use your findings to validate or enrich the final response.

4. Ensure confidence. That is, only make a final decision when you are confident that the retrieved information fully addresses the user's query.

1101

## F.5 Task Prompt

The task prompt defines the concrete input for the current user question, which should at least include: 1) the user input question, 2) the answer format, and 3) the database or vectorstore schema (if needed). Following Nan et al. (2023), we use the code representation format for database schema and incorporate schema descriptions to enhance the schema linking. As for the vectorstore schema, we introduce 1) all available collections, 2) fields for each stored data entry, 3) encodable (table, column) pairs from the corresponding DuckDB, and 4) valid operators to use in the filter condition during vector search. We only present one input case for NeuSym-RAG, while the task prompt for other methods can be easily inferred depending on whether the database or vectorstore will be integrated as the backend environment.

1102

## Task Prompt for NeuSym-RAG

Remember that, for each question, you only have 20 interaction turns at most. Now, let's start!

**[Question]:** What are the main questions that this paper tries to resolve or answer?"

**[Answer Format]:** Your answer should be a Python list of text strings, with each element being one critical problem that this paper analyzes, e.g., ["question 1", "question 2", ...].

**[Database Schema]:** The database schema for "ai\_research" is as follows:

```
/* database ai_research: This database contains information about AI research papers. Each PDF file is
   represented or parsed via different views, e.g., pages, sections, figures, tables, and references.
   We also extract the concrete content inside each concrete element via OCR. */
/* table metadata: This table stores metadata about each paper. */
CREATE TABLE IF NOT EXISTS metadata (
    paper_id UUID, -- A unique identifier for this paper.
    title VARCHAR, -- The title of this paper.
    abstract VARCHAR, -- The abstract of this paper.
    pub_year INTEGER, -- The year when this paper was published.
    ... [more columns with their data types and descriptions, omitted] ...
    PRIMARY KEY (paper_id)
);
... [the remaining tables and columns using the CREATE statement] ...
```

**[Vectorstore Schema]:** The vectorstore schema for ai\_research is as follows. You can try collections with different encoding models or modalities:

```
[
  {
    "collection_name": "text_bm25_en",
    "description": "This collection is used to store the sparse embeddings generated by the BM25 model for all encodable text content in another relational database. The semantic search is based on field `vector` with metric inner-product (IP).",
    "fields": [
      {"name": "vector", "dtype": "SPARSE_FLOAT_VECTOR", "desc": "attained by BM25 model"},
      {"name": "text", "dtype": "VARCHAR", "desc": "cell value from the database"},
      {"name": "pdf_id", "dtype": "VARCHAR", "desc": "unique id of the PDF file"},
      {"name": "page_number", "dtype": "INT16", "desc": "source page of the `text` field"},
      {"name": "table_name", "dtype": "VARCHAR", "desc": "source table of `text` field"},
      {"name": "column_name", "dtype": "VARCHAR", "desc": "source column of `text` field"},
      {"name": "primary_key", "dtype": "VARCHAR", "desc": "primary key value for the row
          that contains the `text` field in the relational database"}
    ]
  },
  {
    "collection_name": "text_sentence_transformers_all_minilm_l6_v2",
    "description": "This collection is used to store the embeddings generated by the model MinILM-L6-v2 model for all encodable text content in another relational database. The semantic search is based on field `vector` with metric COSINE.",
    "fields": "The fields of this collection are the same as those in `text_bm25_en`."
  },
  ... [other collections with their fields] ...
]
```

Here are all encodable (table\_name, column\_name) tuples from the corresponding DuckDB database, where the encoded vector entries are sourced. Different columns together provide multiple perspectives for vector search.

```
[("metadata", "title"), ("metadata", "abstract"), ... [more encodable (table, column) pairs] ...]
```

Here are the operators that you can use in the filter parameter for RetrieveFromVectorstore action:

```
[
  {
    "symbol": "and",
    "example": "expr1 and expr2",
    "description": "True if both expr1 and expr2 are true."
  },
  {
    "symbol": "+",
    "example": "a + b",
    "description": "Add the two operands."
  },
  ... [more operators with detailed description] ...
]
```

1112

## G Summary of Changes

1113 We make the following modifications in the new version:

- 1114 1. Annotate and increase the size of our self-annotated dataset AIRQA-REAL from 100 to 553. It also  
1115 extends from the original single PDF based samples to more diverse “multi-doc analysis” and “paper  
1116 retrieval” settings. The detailed statistics are elaborated in App. A;
- 1117 2. Convert and experiment on two other existing full-PDF based Q&A datasets, namely M3SciQA  
1118 (Li et al., 2024a) and SciDQA (Singh et al., 2024), to validate the universality of the proposed  
1119 NeuSym-RAG framework. They contain 452 and 2937 test samples respectively;
- 1120 3. Experiment with more state-of-the-art open-source large language models, including Qwen2.5-VL-  
1121 Instruct (Team, 2025), InternLM2.5 (Cai et al., 2024), and DeepSeek-R1 (Guo et al., 2025);
- 1122 4. Conduct extra ablation study in Figure 5 about the variation of different model sizes to verify whether  
1123 the performance scales well with parameters and explore the feasibility of utilizing smaller LLMs;
- 1124 5. Implement more baselines, including trivial inputs and structured methods in Table 3. These baselines  
1125 include: a) Question only; b) Title+Abstract; c) Full-text w/. cutoff; d) GraphRAG (Edge et al., 2024);  
1126 e) HybridRAG (Sarmah et al., 2024). Taking into account the previously implemented baselines, we  
1127 have implemented a total of 11 methods.
- 1128 6. Refine the logic of writing and re-draw the main framework Figure 2 and insert a new Figure 4 to  
1129 visualize different baseline methods.
- 1130 7. Adding more details in the Appendix to elaborate how we convert and utilize two existing Q&A  
1131 datasets (M3SciQA and SciDQA) and how we implement the various baselines.

1132 Furthermore, our code and data will be publicly available and here is an anonymous link to the repository:  
1133 <https://anonymous.4open.science/r/neu-sym-rag/>.