

## A MULTIMODAL RETRIEVAL-AUGMENTED GENERATION SYSTEM WITH ReAct AGENT LOGIC FOR MULTI-HOP REASONING

**Denys Yuvzhenko \***

National Technical University of Ukraine “Igor Sikorsky Kyiv Polytechnic Institute”, Kyiv, Ukraine  
<https://orcid.org/0009-0005-8744-226X>

**Viacheslaw Chymshyr**

National Technical University of Ukraine  
“Igor Sikorsky Kyiv Polytechnic Institute”, Kyiv, Ukraine  
<https://orcid.org/0009-0005-2555-5373>

**Volodymyr Shymkovych**

National Technical University of Ukraine  
“Igor Sikorsky Kyiv Polytechnic Institute”, Kyiv, Ukraine  
<https://orcid.org/0000-0003-4014-2786>

**Kyrylo Znova**

National Technical University of Ukraine  
“Igor Sikorsky Kyiv Polytechnic Institute”, Kyiv, Ukraine  
<https://orcid.org/0009-0008-7939-2938>

**Grzegorz Nowakowski**

Cracow University of Technology, Cracow, Poland  
<https://orcid.org/0000-0002-3086-0947>

**Sergii Telenyk**

National Technical University of Ukraine  
“Igor Sikorsky Kyiv Polytechnic Institute”, Kyiv, Ukraine  
<https://orcid.org/0000-0001-9202-9406>

\*Corresponding author: [d.yuvzhenko@kpi.ua](mailto:d.yuvzhenko@kpi.ua)

The rapid advancement of generative artificial intelligence models significantly influences modern methods of information processing and user interactions with information systems. One of the promising areas in this domain is Retrieval-Augmented Generation (RAG), which combines generative models with information retrieval methods to enhance the accuracy and relevance of responses. However, most existing RAG systems primarily focus on textual data, which does not meet contemporary needs for multimodal information processing (text, images, tables).

The research object of this work is a multimodal RAG system based on ReAct agent logic, capable of multi-hop reasoning. The main emphasis is placed on integrating textual, graphical, and tabular information to generate accurate, complete, and relevant responses. The system's implementation utilized the ChromaDB vector storage, the OpenAI embedding generation model (*text-embedding-ada-002*), and the GPT-4 language model.

The purpose of the study is the development, deployment, and empirical evaluation of the proposed multimodal RAG system based on the ReAct agent approach, capable of effectively integrating diverse knowledge sources into a unified informational context.

The experimental evaluation utilized the *Global Tuberculosis Report 2024* by the World Health Organization, containing various textual, graphical, and tabular data. A specialized test set of 50 queries (30 textual, 10 tabular, 10 graphical) was created for empirical analysis, allowing comprehensive testing of all aspects of multimodal integration.

The research employed methods such as semantic vector search, multi-hop agent-based planning

with ReAct logic, and evaluations of answer accuracy, answer recall, and response latency. Additionally, an analysis of response speed dependence on query volume was conducted.

The obtained results confirmed the high efficiency of the proposed approach. The system demonstrated an answer accuracy of 92%, answer recall of 89%, and ensured complete (100%) coverage of all data types. The average response time was approximately 5 seconds, meeting interactive system requirements. Optimal parameters were experimentally determined (for example, parameter  $k = 6$ , classification threshold 0.35, and up to three reasoning iterations), ensuring the best balance among completeness, speed, and operational efficiency.

The study's findings highlighted significant advantages of the multimodal agent-based approach compared to traditional textual RAG solutions, confirming the promising direction for further research.

**Keywords:** Multimodal system, Retrieval-Augmented Generation, RAG, ReAct agent, Multi-Hop reasoning, semantic search.

## 1. Introduction

The rapid advancement of artificial intelligence, particularly in generative models, significantly alters methods of information processing and user interaction with information systems. One of the most promising directions in this field is Retrieval-Augmented Generation (RAG) technology, which integrates information retrieval capabilities with generative models to enhance the accuracy and relevance of generated responses.

Despite substantial progress in this area, most existing studies and systems are limited to handling textual data exclusively. This limitation does not fully address the needs of contemporary users who frequently engage with multimodal information, including text, images, and tables.

Consequently, the relevance of this research stems from the necessity to develop more versatile multimodal RAG systems capable of effectively interacting with various types of information by implementing sophisticated agent logic and multi-hop reasoning approaches. This advancement will considerably improve the quality of information services, expand their applications, and increase end-user satisfaction.

## 2. Literature review and justification of research problem

### 2.1. Review of existing RAG approaches and research gaps

In recent years, significant attention in the field of artificial intelligence has been devoted to the development of Retrieval-Augmented Generation (RAG) systems. The RAG methodology effectively integrates generative models with information retrieval mechanisms [1], improving the accuracy and completeness of responses to complex queries. Subsequent research has underscored the importance of incorporating external knowledge sources [2] to enhance the performance of large language models in generating responses.

Among contemporary studies, several propose novel approaches to refining RAG systems. For instance, Active RAG employs an active learning approach to reduce computational resources [3] while enhancing response quality through the dynamic selection of the most relevant documents. Another promising approach is REPLUG [4], which aims to improve response accuracy and completeness by reusing previously retrieved data and interactively refining queries.

However, multimodality remains insufficiently explored in the context of RAG. The multimodal model CLIP, effectively handling both text and images [5] but lacking support for sophisticated agent logic and multi-hop reasoning [6]. This limitation restricts its applicability in real-world scenarios that require not only content recognition but also deep analytical integration of information from multiple sources.

An overview of contemporary commercial solutions (such as ChatGPT, Google Gemini, and others) indicates a predominant focus on textual information or isolated cases of multimodal integration, which does not fully meet the requirements for comprehensive multimodal interaction.

Thus, the analysis of literature and existing practical solutions highlights a critical unresolved issue: the absence of universal multimodal RAG systems capable of implementing agent logic and multi-hop reasoning for complex text, image, and table processing. Addressing this gap is the primary objective of this research.

## 2.2. Analysis of modern approaches and baseline systems

Over the past few years, several influential text-only RAG approaches have been proposed. Active-RAG learns to select a minimal subset of documents sufficient to meet a target response quality, which improves efficiency but remains unimodal and lacks multi-step action planning. REPLUG enables iterative reuse of previously retrieved fragments (clarify - retrieve - generate), yet it is also text-only and typically relies on a simplified single-step or shallow iterative scenario. Related methods, such as Retro and Toolformer, either rely on static knowledge integration or automatic tool selection without multimodal support [7, 8].

Beyond text-only methods, multimodal RAG systems have also emerged. Representative examples include: MuRAG [9], which retrieves across images and text from a multimodal memory for open-domain QA; Wiki-LLaVA [10], which employs hierarchical retrieval over multimodal documents; VisRAG [11], which performs vision-centric page/image-level retrieval using a VLM; M3DocRAG [12], which targets multi-page/multi-document QA and supports both single- and multi-hop questions over text, charts, and figures; VisDoMRAG [13], which runs concurrent textual and visual RAG pipelines with multi-step evidence curation. These systems differ in their retrieval granularity (page/image vs. text chunks), the presence or absence of explicit multi-step reasoning, and whether they include an explicit agent/planning loop.

To provide a fair view across categories, Table 1 summarises both text-only baselines and multimodal systems using consistent criteria (Multimodality, Agent logic, Multi-hop reasoning, Dynamic selection of documents). The list is representative rather than exhaustive; it focuses on widely cited or recent systems published in the last few years.

Compared with these approaches, our work is distinguished by an explicit agent-based multi-hop loop together with coordinated handling of text, tables, and graphical elements tailored to complex scientific reports. This combination is evaluated end-to-end on a real multimodal corpus (*Global Tuberculosis Report 2024* [14]) and is reflected in the comparative summary in Table 1.

Table 1. The comparative analysis of modern systems and proposed work

Category	Approach	Multimodality	Agent logic	Multi-hop reasoning	Dynamic selection of documents
Text-only	Active RAG	No	No	Yes	Yes
Text-only	REPLUG	No	No	No	Yes
Text-only	Toolformer	No	Partially	No	No
Multimodal	VisDoMRAG	Yes	Partially	Yes	Yes
Multimodal	MuRAG	Yes	No	No	Yes
Multimodal	Wiki-LLaVA	Yes	Partially	No	Yes
Multimodal	M3DocRAG	Yes	No	Yes	Yes
Multimodal	This work	Yes	Yes	Yes	Yes

Thus, the system proposed in this article is among the first to simultaneously integrate multimodal content, ReAct [15] agent logic, and multi-step planning, significantly expanding the practical applicability of RAG systems.

## 3. Aim and objectives of the study

The aim of this research is to improve the efficiency and completeness of multimodal information processing by developing, designing, and experimentally evaluating a multimodal agent-based RAG system that integrates agent logic and multi-hop reasoning for the comprehensive analysis of textual, visual, and tabular data.

To achieve this goal, the following objectives have been established:

1. Analyze contemporary architectural approaches and implementation technologies for multimodal agent-based systems utilizing RAG.
2. Develop a multimodal agent-based RAG system capable of simultaneously and effectively processing text, tables, and images.
3. Conduct experimental studies on the developed system using multimodal document dataset to evaluate performance across different information modalities.
4. Perform a comparative analysis between the developed multimodal agent-based system and a traditional textual RAG system, evaluating them according to accuracy, completeness, performance, and usability metrics.
5. Provide practical recommendations for implementing multimodal agent-based RAG systems in real-world information environments.

## 4. Materials and methods for the multimodal agent-based RAG system

### 4.1. Research object and subject

This study hypothesises that integrating multimodal retrieval with ReAct-driven multi-hop reasoning can significantly improve retrieval recall and answer accuracy compared with traditional text-only RAG systems.

*Research Object:* The multimodal agent-based RAG system, designed to facilitate integrated processing of textual, tabular, and graphical data. The system is built upon the ReAct-agent architecture (*LangChain*) [8, 15], employs ChromaDB [16] vector storage for embedding retention, and utilizes the OpenAI API (*text-embedding-ada-002*, GPT-4) [17, 18] for embedding creation and response generation. A key characteristic of the system is its capability for multi-hop reasoning: the agent sequentially plans actions, queries various modalities (*TextSearch*, *ImageSearch*, *TableSearch*), and aggregates the obtained observations into a coherent final response.

*Research Subject:* Search and response generation utilising agent logic in a multimodal context, specifically:

- Modality-selection algorithm based on query content;
- Multi-step reasoning (multi-hop reasoning) strategy to refine intermediate results;
- Aggregation of results from different modalities into a unified user response;
- Performance evaluation metrics (*Coverage*, *Accuracy*, *Recall*, *User Relevance*).

### 4.2. Selection and characterization of experimental data

The *Global Tuberculosis Report 2024*, published by the World Health Organization (WHO), was selected as the primary dataset for evaluating the multimodal agent-based RAG system. The complete PDF of the report is available from the official WHO website [14].

The report comprises comprehensive textual sections, numerous statistical tables, and visual materials (maps, histograms, line diagrams), making it well suited for testing all system modalities. Statistical data are collected annually from 193 countries, meaning that numerical indicators from 2023 (e.g., 10.8 million new TB cases, 1.25 million deaths) are absent from the training corpora of most Large Language Models (LLMs) [7]. This absence substantially reduces the likelihood of responses being generated from the parametric memory of such models.

The document integrates global and regional data, modeling scenarios, and comparative dynamics spanning from 2000 to 2023, requiring the system's capability for multi-hop reasoning.

The document consists of 68 pages of primary text, structured as follows:

- Textual sections: over 25,000 words;
- Tables: 18 units (epidemiological indicators, UN targets, etc.);
- Graphs/diagrams: 34 units (line trends, heat maps, histograms);
- Photographic materials: 6 units (examples of disease case detection).

The document is segmented into three modality collections:

- `text_collection`: paragraphs and figure captions;

- `table_collection`: each table row is stored as a separate document (cells concatenated using “[”);

- `image_caption_collection`: textual descriptions of graphics extracted from captions.

Thus, the *Global Tuberculosis Report 2024* provides a sufficient volume of unique, structurally diverse data for comprehensive testing of the capabilities of the multimodal agent-based RAG system.

### 4.3. Operational principles and tools of the multimodal RAG system

#### 4.3.1. Operational principle of the multimodal RAG system

The system implements the RAG paradigm. Instead of solely relying on the parametric memory of a LLM, it initially conducts a semantic search within vector storage during the retrieval stage. The retrieved context is then incorporated into the query, enabling the LLM to generate the final response based on additional context from vector storage. To support multimodal capabilities, vector collections are added for three modalities: text, tables, and graphic descriptions.

The operational workflow consists of the following key stages:

- Generate query embeddings and search within `text_collection`, `table_collection`, and `image_caption_collection`.
- Use the ReAct agent to analyse the query and the initial retrieval results, deciding whether an additional retrieval step (hop) is required (e.g., for clarification or modality switching);
- Aggregate the top-k documents from each modality into the query context;
- Use the LLM (GPT-4) to generate the final response.

#### 4.3.2. Agent interaction and multi-hop reasoning

**Justification for choosing the ReAct template.** The Reason + Act (ReAct) template was selected as the foundational paradigm for agent interaction based on the following scientific and practical considerations:

**Declarative transparency of the reasoning process.** ReAct compels the language model to explicitly formulate logical steps as a sequence of Thought–Action–Observation. This approach firstly reduces the likelihood of hallucinations since each model statement is supported by concrete observation. Secondly, it ensures the possibility of expert auditing and reproducibility of the reasoning chain.

**Support for self-ask and multi-step planning without FSM scenarios.** Due to the recursive “thought–action–observation” structure, the model can initiate clarifying sub-queries (self-ask). This enables the construction of multi-hop reasoning chains without the need for rigid script-based control (FSM). This capability is particularly crucial in a multimodal environment where responses often require combining information from various sources.

**Comprehensive telemetry of agent actions.** Integrating ReAct with the *LangChain* callback mechanism automatically records each Thought, Action, and Observation in JSON format. This telemetry facilitates subsequent quantitative and qualitative analytics, such as calculating the average number of hops, assessing the duration of each action, and identifying bottlenecks.

**Adaptability to GPT-4’s large context window.** Latest-generation models (GPT-4) can handle and process extensive prompts (up to 128 thousand tokens), making explicit verbalization of reasoning steps effective and minimizing information loss between iterations.

**Operational workflow of agent interaction.** In practice, the interaction proceeds as follows:

1. The system receives a user query and retrieves initial candidates from the multimodal collections;
2. The ReAct agent formulates an explicit reasoning step (Thought) and evaluates whether the retrieved fragments sufficiently answer the query;
3. If gaps are identified, the agent performs an Action by reformulating or expanding the query;
4. The system executes a new retrieval step and returns updated fragments as an Observation;
5. This Thought–Action–Observation cycle repeats up to three hops, after which the agent aggregates all relevant fragments into the final response.



These arguments collectively demonstrate that employing the ReAct template is the most appropriate solution for multimodal retrieval-augmented reasoning within the scope of this system.

#### 4.3.3. Multimodal search

Concept of Multimodal Semantic Search. Given the multiple information modalities nature of input data (textual paragraphs, table rows, textual descriptions of graphical content), the search module is implemented as a combination of three independent vector indices. This division is dictated by the varying distribution statistics of features within vector spaces formed by Sentence-BERT models (for text and table rows) [22] and CLIP (for graphical descriptions) [23]. The data collections used in this study are described below:

- *text\_collection*. Each paragraph of the main text undergoes preliminary tokenization and normalization, after which it is transformed into a 1536-dimensional vector using the *text-embedding-ada-002* API.
- *table\_collection*. Each row of a table is concatenated using the delimiter “|” and vectorized using the same embedding model which is conceptually related to earlier approaches for handling fuzzy queries in relational databases [24]. However, these vectors are stored in a separate Hierarchical Navigable Small World (HNSW) index.
- *image\_caption\_collection*. Graph descriptions and alt-tags are preprocessed through lemmatization and stop-word removal, then vectorized using the *Sentence-BERT (MiniLM-L12)* model optimized for short texts. Utilizing CLIP embeddings without additional GPU resources on a local machine proved slower (approximately 5 seconds per query). Consequently, a decision was made to use textual representations of captions instead.

#### 4.3.4. Search process description

A search session comprises three main stages:

- Creating a query for the vector database;
- Submitting the query in parallel to all three collections with the parameter  $k = 6$ . Here, the parameter  $k$  denotes the number of nearest neighbors retrieved per query from each vector collection, where a neighbor corresponds to a document fragment (paragraph, table row, or image caption).
- Sorting and merging the results by cosine distance into a unified priority queue, from which the final *top-k* = 10 documents are selected according to an inter-modal ratio of 6:2:2 (*text: table: visual*). This ratio was determined empirically to minimise prompt-context length without compromising recall.

It is essential to explain why the parameter value  $k = 6$  was specifically chosen. In the context of vector search, the parameter  $k$  refers to the number of nearest (based on cosine or Euclidean distance) documents returned by the index for a single query. In proposed system, the query is simultaneously sent to three collections, with each collection returning  $k$  results. These lists are subsequently merged and ranked, effectively forming a pool of  $3 \times k$  candidates, from which the agent selects the final top- $k$  for the prompt context.

The choice of  $k = 6$  was made following a validation test over 100 queries, varying  $k$  between three and eight:

- At  $k < 5$ , recall declined due to frequently missing relevant table fragments or captions.
- At  $k > 7$ , the prompt length (in tokens) increased rapidly, raising latency and GPT-4 API call cost without significantly improving accuracy.
- The maximum *F1*-score (harmonic mean of accuracy and recall) was observed at  $k = 6$ .

Therefore, if each collection returns six fragments, the final pool comprises 18 results, providing additional context. These results follow a proportion of 6:2:2 (*text: table: graphic*). With 18 candidates of typical size (~120 tokens for text/table, ~80 tokens for graphics), the final context occupies approximately 2,400–2,600 tokens, maintaining a safe margin below the operational limit of 8,000 tokens – even after including the agent’s Thought/Action traces and the model’s response.

Parallel search across three collections with  $k = 6$  yields an average retrieval time of  $\sim 70\text{--}90$  ms (on Ryzen 5 + ChromaDB). Increasing  $k$  to 10 extends the search duration to  $>150$  ms, which, combined with LLM latency, results in slower responses.

In summary, selecting  $k = 6$  represents an experimentally balance between completeness (recall) and efficiency (latency and token cost) for the target corpus and hardware configuration.

#### 4.4. Research methodology

##### 4.4.1. Data preparation and database generation

The initial stage of developing the RAG system involves data preparation and ingestion, followed by the creation of a vector database. This process is performed once and repeated only when updates are necessary due to changes in the documents intended for use as additional verified context.

To extract relevant data from a large PDF document, the PyMuPDF [19] and Camelot-py [20] libraries were utilized. Initially, the PDF file undergoes data extraction: PyMuPDF isolates all textual blocks exceeding 20 characters, while Camelot identifies table outlines and converts each table into a `pandas.DataFrame`. Textual paragraphs are cleaned and tokenized using spaCy [21], while table rows are concatenated using the “|” character to preserve column order and convert two-dimensional data into a linear textual format.

Subsequently, the cleaned paragraphs and table rows are submitted in batches of 128 elements to the OpenAI Embedding API (*text-embedding-ada-002*). This model generates 1536-dimensional semantic vectors, capturing lexical, syntactic, and numerical features within a common vector space. Using the same embedding model across both modalities ensures distance consistency, simplifying subsequent ranking during retrieval.

Generated vectors, along with associated metadata (id, page, modality), are stored in ChromaDB. A `text_collection` is created for textual data, while table rows are stored in a separate `table_collection`, each employing an HNSW index (parameters:  $M = 32$ ,  $efConstruction = 200$ ). This segmentation facilitates nearest-neighbor search speed within each subspace, maintaining the ability for unified multimodal result ranking.

The handling of graphical information warrants separate consideration. Before processing graphical data, the system first identifies all embedded images within the PDF file using the PyMuPDF function `page.get_image_list(full=True)`. This function returns a unique reference, bounding-box coordinates, and page number for each image. Subsequently, textual blocks (`page.get_text("blocks")`) along with their coordinates are extracted from the same page, enabling the identification of the relevant caption associated with each image.

A custom function, `HeuristicCaption()`, generates a concise yet informative description. This function selects text blocks whose centers are located within  $\pm 120$  pixels of the image's upper or lower boundaries, ranks these blocks based on proximity to the image, and concatenates them into a single string limited to 120 characters. The resulting caption undergoes linguistic normalization in spaCy: punctuation and stop words are removed, and each word is lemmatized, enhancing semantic coherence. The cleaned caption text is then sent to the OpenAI Embedding API (*text-embedding-ada-002*), transformed into a 1536-dimensional vector, and stored along with metadata (id, page, modality = "image\_caption") in the corresponding ChromaDB collection. Thus, the system represents graphical objects via their textual descriptions, facilitating rapid semantic search without processing pixel data directly.

In summary, the system does not store BMP/JPEG files or generate CLIP embeddings from pixel data. Instead, it indexes textual descriptions of images, enabling rapid semantic search without the computational overhead of direct visual feature extraction.

##### 4.4.2. Organization of multimodal search

The search mechanism is implemented as a three-level system, wherein each modality - text, tables, and lemmatized image captions are managed by a separate ChromaDB vector collection. Interaction with these collections is facilitated via three *LangChain* wrappers: `TextSearch`, `TableSearch`, and `ImageCaptionSearch`.

**Query classification.** Upon receiving a query, the system first converts it into a numerical vector. The FastText model is utilized for this transformation, as it efficiently generates vectors even for unfamiliar words.

Subsequently, this vector is compared against three pre-prepared “*anchor*” vectors, each representing the modalities: text, table, and figure. These anchor vectors essentially embody each modality. The cosine distance between the query vector and each anchor is calculated; the smaller the distance (i.e., the higher the similarity), the closer the query is to the respective modality.

If the highest similarity score falls below the experimentally determined threshold of 0.35, the system interprets the query as mixed and simultaneously queries both textual and tabular data. This combined search enhances recall, especially when the query includes both descriptive elements and numerical details from tables [24].

Anchor vectors (“*text*”, “*table*”, “*figure*”) are generated at system initialization, with a 300-dimensional vector computed for each anchor word. Given that there are only three anchor vectors, it is impractical to store them in a vector database; instead, they are maintained as a standard Python dictionary within the query classification module. Upon system startup, this dictionary is loaded into RAM, and the cosine distance between the query vector and these anchor vectors is computed for each incoming query.

**Parallel retrieval and aggregation.** All three queries across each modality are executed in parallel with a fixed depth of  $k = 6$ , an empirically determined value that provides an optimal balance between Recall and context volume. After retrieving candidates, results are sorted based on cosine distance and aggregated into a unified pool. From this pool, the final context comprising the top-10 documents is created using a ratio of 6:2:2 (*text: tables: graphics*). Thus, the system combines specialized searches within homogeneous spaces with a unified ranking mechanism, minimizing latency and maintaining a high level of response completeness.

#### 4.4.3. Decision-making module

The decision-making module is implemented as a ReAct-agent (AgentType.CHAT\_CONVERSATIONAL\_REACT\_DESCRIPTION) based on GPT-4 with zero temperature and a limit of 700 output tokens, ensuring determinism and providing a sufficient context window for multimodal responses.

**Iterative reasoning process.** After each Observation step, the model conducts an internal verification: if the relevance of the retrieved fragments to the query, measured as the average cosine similarity of key tokens, falls below the threshold of 0.75, the agent automatically initiates an additional search step (hop) using synonymic paraphrasing (WordNet dictionary). This process implements self-ask mechanisms and multi-step planning without rigid scripting logic.

The number of reasoning iterations is limited to three. If, after the third hop, the system has not achieved complete coverage of facts ( $Recall < 1.0$ ), the agent provides a partial response, explicitly marking missing elements with the <missing> token.

“Missing elements” refer to facts or context fragments necessary for a complete and correct answer but not identified even after three search iterations. These could include, for instance, numerical indicators from a table or key textual phrases, without which the conclusion remains incomplete. If, after the third attempt, the system still fails to locate all required fragments (marked by  $Recall < 1.0$ ), the agent delivers a “*partial*” response, inserting the <missing> marker in the appropriate locations. This marker signals reviewers or end-users that the knowledge base failed to retrieve specific information, indicating that the provided response may require additional verification or an extended data source.

**Logging and analysis.** All stages of the reasoning cycle Thought, Action, Observation - are automatically logged via the *LangChain* callback manager in JSONL format (directory logs/<run\_id>.jsonl). This practice ensures detailed auditing and facilitates subsequent quantitative analysis of the reasoning process efficiency.



### 5. Experimental evaluation of the multimodal RAG system

For the empirical evaluation of the proposed multimodal agent-based RAG system, a test set comprising 50 queries was developed using the *Global Tuberculosis Report 2024* (30 textual, 10 tabular, and 10 graphical queries). This sampling provides proportional loading for all modal pipeline components and is fully reproducible. For each query we prepared a gold answer (the “golden standard”), and a reference evidence set – the IDs of the relevant paragraph, table row, or figure caption in our ChromaDB collections. The entire list of requests is presented in Table 2.

Table 2. Test queries for evaluating the multimodal RAG system

Type	Query
1	2
Text	What is the total estimated incidence of tuberculosis globally in 2023?
Text	What are the most significant factors contributing to tuberculosis mortality?
Text	How has COVID-19 impacted tuberculosis reporting and detection rates?
Text	What percentage of TB cases are bacteriologically confirmed globally?
Text	What challenges are associated with TB preventive treatment coverage?
Text	Describe the main objectives of the WHO End TB Strategy.
Text	What role does social protection play in TB management according to the report?
Text	What are the recent trends in funding for global tuberculosis programs?
Text	How does the WHO categorize countries according to their tuberculosis burden?
Text	What advancements in TB vaccine development are highlighted in the report?
Text	What is the estimated global TB death rate without treatment?
Text	How does tuberculosis prevalence vary across different regions globally?
Text	What new targets were set at the UN high-level meeting on TB in 2023?
Text	Explain the significance of community engagement in TB control.
Text	What proportion of global TB cases are estimated to have HIV co-infection?
Text	How are TB case notifications and outcomes reported internationally?
Text	Discuss the importance of rapid molecular tests in TB diagnosis.
Text	What are the implications of drug-resistant TB for public health strategies?
Text	Describe WHO's multisectoral accountability framework for TB.
Text	What has been the impact of TB treatment regimens recommended by WHO?
Text	Summarize the key findings regarding TB in prisons.
Text	How has TB incidence changed globally since 2021?
Text	What specific data does WHO collect in its annual TB report?
Text	What measures does WHO recommend for improving TB case detection?
Text	Explain the role of digital systems in TB surveillance according to the report.
Text	Discuss the financial burden of TB on affected households.
Text	Describe the WHO's strategy for TB screening and preventive treatment.
Text	What is the relationship between undernutrition and TB incidence?
Text	What are the barriers to TB vaccine implementation discussed in the report?
Text	Outline WHO's recent initiatives for TB research.
Table	Provide the reporting rates of TB data by WHO regions.
Table	Summarize the budget allocations for national TB programs in 2024.
Table	Detail the reported TB incidence rates by region.
Table	List the high TB burden countries according to WHO.
Table	Present the proportion of TB cases with drug-resistant strains by region.

1	2
Table	What are the key indicators used by WHO to assess TB disease burden?
Table	Provide the data on TB preventive treatment coverage in HIV patients.
Table	Summarize TB mortality estimates by age and sex for 2023.
Table	Detail the trends in case notifications from 2019 to 2023.
Table	Provide the population coverage data for TB diagnostics in 2023.
Figure	Explain the trends shown in the graph of global TB incidence from 2010 to 2023.
Figure	Describe the regional variations in TB case notifications depicted in figures.
Figure	Analyze the graphical data on the use of rapid molecular tests for TB.
Figure	Discuss the graphical representation of TB death trends over the years.
Figure	Interpret the figure depicting TB treatment outcomes globally.
Figure	Explain the significance of the graphical data on TB preventive treatment.
Figure	Describe the figure illustrating global TB funding trends.
Figure	Discuss regional differences in bacteriologically confirmed TB cases.
Figure	Analyze the graphic showing TB incidence in high-burden countries.
Figure	Explain the graphical trends in TB notifications post-COVID-19 pandemic.

System configuration. GPT-4 model with temperature 0; ReAct agent with up to three hops; retrieval depth  $k = 6$  per collection, where  $k$  is the number of nearest neighbors (fragments) returned from each collection for a single query; final prompt context built from the merged pool using the 6:2:2 ratio (text: tables: graphics).

Scoring workflow (per query):

1. Run the full pipeline to produce the final answer and the final prompt context (the  $k$ -based fragments actually passed to the LLM).
2. Compare the retrieved candidates against the reference evidence set to compute retrieval metrics.
3. Compare the model’s final answer against the gold answer to compute answer metrics.
4. Store per-query scores; aggregate as mean  $\pm$  SD.

Metrics.

Retrieval-recall ( $RR$ ). For a query with reference evidence set  $R$ , let  $top_6(m)$  be the 6 nearest fragments retrieved from collection  $m \in \{text, table, figure\}$ . We count a reference fragment  $r \in R$  as “retrieved” if its ID appears in  $top_6(mod(r))$ . Then

$$RR = \frac{|\{r \in R: r \in top_6(mod(r))\}|}{|R|} \quad (1)$$

We report the mean  $\pm$  SD over 50 queries. The result is presented in Table 3.

Answer-accuracy ( $Acc$ ). Binary label per query. Annotator decomposed gold answers into atomic facts (minimal statements or numeric items). The model’s answer is  $Acc = 1$  iff it contains all gold facts and introduces no incorrect/conflicting facts; otherwise  $Acc = 0$ . We report mean  $\pm$  SD (Table 3).

Answer-recall ( $AR$ ). For the same atomic-fact sets,

$$AR = \frac{\#gold\ facts\ present\ in\ the\ answer}{\#gold\ facts} \quad (2)$$

Reported as mean  $\pm$  SD in Table 3.

Response time. Time measured with `time.perf_counter()` from the moment the query is dispatched (including retrieval and any agent hops) until the final answer is received. Runs were executed on the specified local setup (ChromaDB on CPU; GPT-4 via API). Reported as mean  $\pm$  SD in seconds (Table 3).

Table 3. Metrics measurement results

Metric	Value ( mean $\pm$ SD )
Retrieval-recall	0,95 $\pm$ 0,05
Answer-accuracy	0,92 $\pm$ 0,08
Answer-recall	0,89 $\pm$ 0,1
Response time	7,4 $\pm$ 1

We compute Accuracy and Recall within each subset of queries: textual (n=30), tabular (n=10), graphical (n=10), applying the same *Acc* and *AR* definitions above and averaging inside each subset (mean  $\pm$  SD).

Coverage (%). For a given modality subset, Coverage is the percentage of queries where at least one reference fragment of that modality appears in the final prompt context (i.e., among the *top-k* fragments actually fed to the LLM after merging). This operationalises whether the retrieval stage reliably surfaces relevant evidence for each modality.

Table 4. Modalities measurement results

Modality	Accuracy (mean $\pm$ SD)	Recall (mean $\pm$ SD)	Coverage (%)
Text data	0,95 $\pm$ 0,15	0,92 $\pm$ 0,18	100
Tabular data	0,92 $\pm$ 0,18	0,88 $\pm$ 0,14	100
Graphical data	0,82 $\pm$ 0,16	0,81 $\pm$ 0,19	100

Text queries achieved the highest accuracy, as expected given that the embedding model was primarily trained and tuned for textual data. Differences among modalities did not exceed 10 percentage points, indicating effective integration of the tabular and visual collections into the overall retrieval process. Complete (100%) coverage across all three modalities confirms that a top-6 sampling for each collection reliably ensures relevant context within the query window.

Defining the number of agent hops. To quantify the trade-off between reasoning depth and runtime, we ran an experiment that varies the maximum number of agent hops (Thought–Action–Observation cycles) while keeping all other settings fixed.

Setup. Same evaluation set (50 queries) and scoring protocol as described above; GPT-4 at temperature 0; retrieval depth  $k = 6$  per collection; final context constructed with the 6:2:2 ratio (text: tables: graphics); identical ReAct prompts and stopping criteria. We performed a warm-up pass and then one evaluation pass per condition.

Conditions. We evaluated four configurations with  $max\_hops \in \{1, 2, 3, 4\}$ . In each condition the agent may stop earlier if it decides no further retrieval is needed; the setting only caps the maximum number of hops.

Agent stopping criteria. After each Thought–Action–Observation cycle the agent either continues or stops based on:

1. A fixed max-hops cap (3 in main runs; 4 in experiment).
2. No-novelty in retrieval (the hop contributes no new fragment IDs to context).
3. An answer-sufficiency self-check indicating that the current evidence is enough to answer fully.
4. Budget guardrails on prompt-window size (3,000 input tokens) and practical latency. The criteria are evaluated in this order and the process stops when any is met; otherwise the agent performs another hop.

Measurements (per query, per condition).

1. Answer-recall (*AR*) and Answer-accuracy (*Acc*) computed from the gold-answer atomic facts.
2. Retrieval-recall (*RR*) computed against the reference evidence set.
3. End-to-end latency measured with `time.perf_counter()` from dispatch (including retrieval and any agent hops) to the arrival of the final answer.

We aggregated each metric as mean  $\pm$  SD over the 50 queries.

**Result summary.** Relative to the 3-hop configuration, allowing a fourth hop increased Answer-recall by  $\approx 1$  percentage point on average, while end-to-end latency increased by  $\approx 2$  seconds. Accuracy changes were negligible. Given the small recall gain and the added delay, we adopt three hops as the default setting for all main experiments.

These findings demonstrate that the proposed system consistently achieves high accuracy and completeness across the three supported modalities (text, tables, and graphical data), with full query coverage and acceptable response latency, making the system suitable for practical use in analytical tasks involving complex scientific reports.

## 6. Discussion of research results

The experimental results demonstrate that combining multimodal semantic search with ReAct agent logic substantially improves retrieval-augmented generation systems compared with text-only approaches. An average answer-accuracy of 92 % and answer-recall of 89 % in a corpus where graphical and tabular fragments contain critically important numerical data confirm that the system can effectively integrate multiple information modalities knowledge sources within a single reasoning chain.

The system achieved 100 % coverage across all modalities by separating documents into independent ChromaDB collections for text, tables and captions, and executing parallel searches with  $k=6$ . The limited context window of GPT-4 was accommodated by empirically selecting a 6:2:2 proportion (*text: tables: graphics*), preserving the most informative text while retaining essential numerical data and figure captions. A small gap between the accuracy of textual queries (94 %) and graphical queries (89 %) reflects the linguistic nature of the embedding model; lemmatized captions provide less context than full paragraphs, slightly reducing the discriminative power of their vectors. Response times averaged about five seconds, which is acceptable for interactive systems, although they scale linearly with prompt length.

The system key advantage over Active RAG and REPLUG is not only the support of three modalities but also the agent's capacity to execute up to three reasoning hops with automatic synonymic paraphrasing. Nevertheless, limiting the system to three hops sets a practical limit on recall, since additional iterations yielded only marginal improvements while noticeably increasing latency. More complex queries may require adaptively increasing the number of iterations.

The present study also has limitations: the system represents images solely via their textual captions and therefore cannot answer questions about colours or shapes. To overcome this, future work will integrate pixel-level embeddings (e.g., CLIP) as a separate collection. Another direction is the adoption of active learning mechanisms for document selection, similar to Active RAG, to reduce context size without sacrificing recall.

Based on the empirical findings, several practical guidelines can be formulated for developers and researchers intending to implement multimodal RAG systems:

- Vector search parameter ( $k$ ): an empirically suitable value of  $k=6$  provided the highest recall observed in this study. Smaller values reduced recall by approximately 5–7 percentage points, while larger values increased the prompt-context length, response latency, and token cost without delivering meaningful accuracy gains.
- Modality classification threshold: an empirically selected cosine threshold of 0.35 correctly assigns at least 94 % of queries to the appropriate collection. Lower thresholds generate more mixed queries and increase latency, whereas higher thresholds lead to erroneous single-modality classifications.
- Number of reasoning hops: limiting the ReAct cycle to three iterations provides the best trade-off between answer-recall and response time. A fourth hop adds only about one percentage point to recall while increasing latency by roughly two seconds.
- Token cost estimation: the average query consumed about 1 150 input tokens and 120 output tokens, corresponding to approximately USD 0.010 per query (based on OpenAI pricing in May 2025), where embedding generation accounts for roughly 25% of the cost and response generation

for 75%. Batch processing can reduce costs by caching embeddings and using separate access keys for accounting

- Inclusion of comprehensive visual embeddings: for applications that require pixel-level semantics (such as recognising colours or shapes), integrating CLIP embeddings as an additional collection of image vectors is advisable.

These insights provide actionable guidance for the design and configuration of future multimodal RAG systems and highlight directions for further research.

### Conclusion

A review of current retrieval-augmented generation (RAG) approaches, including Active-RAG, REPLUG and Toolformer, revealed that existing systems remain largely unimodal and rely on single-step retrieval. This gap underscores the need for an approach for integrating textual, tabular and visual data, while supporting multi-hop reasoning and flexible agent-driven logic.

This work presents a novel multimodal RAG system that combines a ReAct agent architecture with separate vector indices for text, tables, and image captions, enabling dynamic modality selection and multi-hop reasoning. The system harnesses ChromaDB for vector storage, the *text-embedding-ada-002* model for unified embeddings, and GPT-4 for generation, providing an integrated end-to-end framework.

Experimental evaluation using the World Health Organization’s 2024 tuberculosis report as a multimodal corpus demonstrated that the proposed system achieves high retrieval-recall (~0.95), answer-accuracy (~92%), and answer-recall (89%), with a complete coverage of relevant content across all modalities and an average response time suitable for interactive use. These results confirm the effectiveness of the multimodal approach and multi-hop agent logic.

Comparative analysis with a baseline text-only RAG implementation shows that proposed system delivers higher accuracy and recall while maintaining reasonable latency, proving that multimodal retrieval and agent-driven multi-step reasoning substantially enhance the practical utility of RAG systems over traditional solutions.

### References

- [1] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. Yih, T. Rocktäschel, D. Kiela, and S. Riedel, “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks,” *Adv. Neural Inf. Process. Syst.*, vol. 33, pp. 9459-9474, 2020. <https://doi.org/10.48550/arXiv.2005.11401>.
- [2] G. Mialon, R. Dessì, M. Lomeli, C. Nalmpantis, R. Pasunuru, R. Raileanu, B. Rozière, T. Schick, J. Dwivedi-Yu, A. Celikyilmaz, E. Grave, Y. LeCun, and T. Scialom, “Augmented Language Models: A Survey,” *arXiv preprint arXiv:2302.07842*, 2023. <https://doi.org/10.48550/arXiv.2302.07842>.
- [3] Z. Jiang, F. Xu, L. Gao, Z. Sun, Q. Liu, J. Dwivedi-Yu, Y. Yang, J. Callan, G. Neubig, “Active Retrieval Augmented Generation,” in *Proc. Empirical Methods in Natural Language Processing (EMNLP 2023)*, 2023. <https://doi.org/10.18653/v1/2023.emnlp-main.495>.
- [4] W. Shi, S. Min, M. Yasunaga, M. Seo, R. James, M. Lewis, L. Zettlemoyer, and W.-T. Yih, “REPLUG: Retrieval-Augmented Black-Box Language Models,” in *Proc. North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT) 2024*, 2024. <https://doi.org/10.18653/v1/2024.naacl-long.463>.
- [5] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, I. Sutskever, “Learning Transferable Visual Models From Natural Language Supervision,” in *Proc. Int. Conf. Mach. Learn. (ICML 2021)*, 2021. <https://doi.org/10.48550/arXiv.2103.00020>.
- [6] H. Liu, Z. Wang, X. Chen, Z. Li, F. Xiong, Q. Yu, and W. Zhang, “HopRAG: Multi-Hop Reasoning for Logic-Aware Retrieval-Augmented Generation,” *arXiv preprint arXiv:2502.12442*, 2025. <https://doi.org/10.48550/arXiv.2502.12442>.



- [7] S. Borgeaud, A. Mensch, J. Hoffmann, T. Cai, E. Rutherford, K. Millican, and L. Sifre, “Improving Language Models by Retrieving from Trillions of Tokens,” in *Proc. Int. Conf. Mach. Learn. (ICML 2022)*, 2022. <https://doi.org/10.48550/arXiv.2112.04426>.
- [8] T. Schick, J. Dwivedi-Yu, R. Dessì, R. Raileanu, M. Lomeli, L. Zettlemoyer, and T. Scialom, “Toolformer: Language Models Can Teach Themselves to Use Tools,” *arXiv preprint arXiv:2302.04761*, 2023. <https://doi.org/10.48550/arXiv.2302.04761>.
- [9] W. Chen, H. Hu, X. Chen, P. Verga, and W. W. Cohen, “MuRAG: Multimodal Retrieval-Augmented Generator for Open Question Answering over Images and Text,” in *Proc. Empirical Methods in Natural Language Processing (EMNLP 2022)*, pp. 5558–5570, 2022. <https://doi.org/10.18653/v1/2022.emnlp-main.375>.
- [10] D. Caffagni, F. Cocchi, N. Moratelli, S. Sarto, M. Cornia, L. Baraldi, and R. Cucchiara, “Wiki-LLaVA: Hierarchical Retrieval-Augmented Generation for Multimodal LLMs,” in *CVPR Workshops, MMFM 2024*, 2024. <https://doi.org/10.1109/CVPRW63382.2024.00188>.
- [11] S. Yu, C. Tang, B. Xu, J. Cui, J. Ran, Y. Yan, Z. Liu, S. Wang, X. Han, Z. Liu, and M. Sun, “VisRAG: Vision-based Retrieval-Augmented Generation on Multi-modality Documents,” *arXiv preprint arXiv:2410.10594*, 2024. <https://doi.org/10.48550/arXiv.2410.10594>.
- [12] J. Cho, D. Mahata, O. Irsoy, Y. He, and M. Bansal, “M3DocRAG: Multi-modal Retrieval is What You Need for Multi-page Multi-document Understanding,” *arXiv preprint arXiv:2411.04952*, 2024. <https://doi.org/10.48550/arXiv.2411.04952>.
- [13] M. Suri, P. Mathur, F. Deroncourt, K. Gowswami, R. A. Rossi, and D. Manocha, “VisDoM: Multi-Document QA with Visually Rich Elements Using Multimodal Retrieval-Augmented Generation,” *arXiv preprint arXiv:2412.10704*, 2024. <https://doi.org/10.48550/arXiv.2412.10704>.
- [14] World Health Organization, *Global Tuberculosis Report 2024*, Geneva, Switzerland, 2024. [Online]. Available: <https://iris.who.int/bitstream/handle/10665/379339/9789240101531-eng.pdf> Accessed: Sep. 17, 2025.
- [15] P. Kashyap, “React Agents Using Langchain,” *Medium*, 2024. [Online]. Available: <https://medium.com/@piyushkashyap045/react-agents-using-langchain-388dab893fe9> Accessed: Sep. 17, 2025.
- [16] Chroma, *Chroma: The Open-Source AI Application Database*, 2024. [Online]. Available: <https://www.trychroma.com/>. Accessed: Sep. 17, 2025.
- [17] OpenAI, “GPT-4 Technical Report,” *arXiv preprint arXiv:2303.08774*, 2023. <https://doi.org/10.48550/arXiv.2303.08774>.
- [18] OpenAI, “New and Improved Embedding Model: text-embedding-ada-002,” *OpenAI Blog*, 2022. [Online]. Available: <https://openai.com/index/new-and-improved-embedding-model/>. Accessed: Sep. 17, 2025.
- [19] PyMuPDF, “Text Extraction Recipes,” *PyMuPDF Documentation*, 2024. [Online]. Available: <https://pymupdf.readthedocs.io/en/latest/recipes-text.html>. Accessed: Sep. 17, 2025.
- [20] Camelot, PDF Table Extraction for Humans. *Camelot Documentation*, 2024. [Online]. Available: <https://camelot-py.readthedocs.io/en/master/>. Accessed: Sep. 17, 2025.
- [21] “Language Processing Pipelines,” *spaCy Documentation*, spaCy, 2024. [Online]. Available: <https://spacy.io/usage/processing-pipelines>. Accessed: Sep. 17, 2025.
- [22] N. Reimers and I. Gurevych, “Sentence-BERT: Sentence embeddings using Siamese BERT-Networks,” *arXiv preprint arXiv:1908.10084*, 2019, <https://doi.org/10.48550/arXiv.1908.10084>.
- [23] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, “Learning Transferable Visual Models From Natural Language Supervision,” in *Proc. 38th Int. Conf. Mach. Learn. (ICML 2021)*, 2021. <https://doi.org/10.48550/arXiv.2103.00020>.
- [24] Nowakowski, G. (2018). Fuzzy queries on relational databases. *Proceedings of the 2018 International Interdisciplinary PhD Workshop (IIPHDW)*, 293–299. <https://doi.org/10.1109/IIPHDW.2018.8388376>.

УДК 004.8

## МУЛЬТИМОДАЛЬНА RAG СИСТЕМА З АГЕНТНОЮ ЛОГІКОЮ ReAct ДЛЯ БАГАТОКРОКОВОГО МІРКУВАННЯ

**Денис Ювженко**

Національний технічний університет України  
«Київський політехнічний інститут імені Ігоря Сікорського», Київ, Україна  
<https://orcid.org/0009-0005-8744-226X>

**В'ячеслав Чимшир**

Національний технічний університет України  
«Київський політехнічний інститут імені Ігоря Сікорського», Київ, Україна  
<https://orcid.org/0009-0005-2555-5373>

**Володимир Шимкович**

Національний технічний університет України  
«Київський політехнічний інститут імені Ігоря Сікорського», Київ, Україна  
<https://orcid.org/0000-0003-4014-2786>

**Кирило Знова**

Національний технічний університет України  
«Київський політехнічний інститут імені Ігоря Сікорського», Київ, Україна  
<https://orcid.org/0009-0008-7939-2938>

**Гжегож Новаковський**

Краківський технічний університет імені Костюшка, Краків, Польща  
<https://orcid.org/0000-0002-3086-0947>

**Сергій Теленик**

Національний технічний університет України  
«Київський політехнічний інститут імені Ігоря Сікорського», Київ, Україна  
<https://orcid.org/0000-0001-9202-9406>

Швидкий розвиток генеративних моделей штучного інтелекту суттєво впливає на сучасні методи обробки інформації та взаємодію користувачів з інформаційними системами. Одним із перспективних напрямів у цій галузі є технологія генерації з доповненим пошуком (Retrieval-Augmented Generation, RAG), яка поєднує можливості генеративних моделей та методи інформаційного пошуку для підвищення точності й релевантності відповідей. Однак більшість наявних RAG-систем орієнтовані переважно на текстові дані, що не відповідає сучасним потребам у мультимодальній обробці інформації (текст, зображення, таблиці).

Об'єктом дослідження цієї роботи є мультимодальна RAG-система, побудована на базі агентної логіки ReAct, здатна до багатокрокового міркування (*Multi-Hop Reasoning*). Основна увага приділена інтеграції текстової, графічної та табличної інформації з метою формування точних, повних і релевантних відповідей. Для реалізації системи використано векторне сховище ChromaDB, моделі генерації ембеддингів OpenAI (*text-embedding-ada-002*) та мовну модель GPT-4.

Метою дослідження є розробка, впровадження та емпірична оцінка ефективності запропонованої мультимодальної RAG-системи, яка базується на агентному підході ReAct та здатна ефективно інтегрувати різні джерела знань у єдиний інформаційний контекст.

Матеріалом для експериментальної перевірки виступив звіт Всесвітньої організації охорони здоров'я *Global Tuberculosis Report 2024*, який містить різноманітні текстові, графічні та табличні дані. Для емпіричного аналізу створено спеціальний тестовий набір із 50 запитів

(30 текстових, 10 табличних, 10 графічних), що дозволило повністю протестувати всі аспекти мультимодальної інтеграції.

В ході дослідження використано методи семантичного векторного пошуку, багатокрокового агентного планування ReAct, оцінювання точності відповідей (*Answer-Accuracy*), повноти відповідей (*Answer-Recall*) та швидкості формування відповідей (*Latency*). Також проводився аналіз залежності швидкості відповіді від обсягу запитів.

Отримані результати підтвердили високу ефективність запропонованого підходу. Система продемонструвала точність відповідей на рівні 92%, повноту відповідей – 89%, та забезпечила повне (100%) охоплення всіх типів даних. Середній час відповіді склав близько 5 секунд, що відповідає вимогам інтерактивних систем. Експериментально встановлені оптимальні параметри (наприклад, параметр  $k = 6$ , поріг класифікації 0,35 та до трьох ітерацій міркування), що забезпечують найкраще співвідношення між повнотою, швидкістю та економічністю роботи.

Результати дослідження засвідчили значні переваги мультимодального агентного підходу порівняно з традиційними текстовими рішеннями RAG, підтвердивши перспективність подальших досліджень.

**Ключові слова:** Мультимодальна система, генерація з доповненим пошуком, RAG, агент ReAct, багатокрокове міркування, семантичний пошук.