

# Building advanced RAG Q&A with Multiple Data Sources Using Langchain: A Multi-Search Agent RAG Application in Ubiquitous Learning

<sup>1</sup>Manel Guettala, <sup>2</sup>Samir Bourekkache, <sup>3</sup>Okba Kazar, <sup>4</sup>Saad Harous

<sup>1</sup>Department of Computer Science,

<sup>1</sup>University of Biskra LINFI Laboratory, Algeria

manel.guettala@univ-biskra.dz, s.bourekkache@univ-biskra.dz,

okba.kazar@sharjah.ac.ae, harous@sharjah.ac.ae

**Abstract**— The integration of Large Language Models (LLMs) in Question-Answering (QA) systems has made significant progress, yet they often fail to generate precise answers for queries beyond their training data and hallucinating. To address this, our study develops an advanced Retrieval-Augmented Generation (RAG) pipeline method using the LangChain framework, featuring a decision-making agent that dynamically selects the most effective tools and data sources for accurate and contextually relevant responses. By incorporating multiple data sources and diverse tools, our system mitigates the limitations of traditional LLMs, enhancing their effectiveness in ubiquitous learning environments. Evaluation through various case studies shows significant improvements in accuracy, relevance, and contextual appropriateness. The multi-search agent RAG system efficiently retrieves and synthesizes information, supporting continuous and context-aware learning and enriching the user experience. This research advances AI-based educational technologies, delivering a powerful solution for information retrieval and synthesis, as well as laying the foundations for intelligent question-and-answer systems of the future.

**Keywords**— *Natural Language Processing (NLP), context-aware Learning, Transfer Learning, Deep Learning, human-computer Interaction (HCI)*

## 1. INTRODUCTION

Recent technological advances and the widespread availability of digital tools have caused big changes in education. It has become very popular to talk about universal learning, a new way of teaching that uses computers to let people learn anytime, anywhere [1]. Learners can access a lot of different kinds of information in a lot of different situations. This model encourages ongoing and flexible learning [2].

Traditional question-and-answer (Q&A) systems have been used a lot in schools to assist students in finding knowledge. Such systems usually use search methods to find relevant texts or parts in a corpus that have already been created, or they use generative models that make responses based on patterns found in the training data [3]. Sometimes more accurate, research-based methods have trouble coming up with logical answers. [4], Answers from generative models can be fluid, even if they aren't exact or useful [5]. Questions-and-answer systems have come a long way since RAG (Retrieval-Augmented Generation) models were introduced [6]. They use relevant papers to write responses, which is a combination of the accuracy of search-based methods and the flexibility of generative models [7].

This study focuses on developing a sophisticated Q&A system using the RAG pipeline, enhanced by multiple search agents to query diverse data sources like academic databases, websites, and proprietary datasets. The main objectives are to create a scalable and reliable system with the Langchain framework, assess its performance in ubiquitous learning environments, and improve the accuracy and relevance of responses. By addressing the limitations of traditional Q&A systems, the proposed approach enhances the learning experience, aligns with the goals of ubiquitous learning, and lays the groundwork for future advancements in AI-driven educational tools.

This paper explores the applicability of Retrieval-Augmented Generation (RAG) in Question-Answering (QA) systems in Ubiquitous Learning Environments (UL) by integrating it with the LangChain framework. Utilizing tools, agents, and vector integration approaches to increase the precision and consistency of information retrieval, demonstrates the synergy between RAG and LangChain. The structure of our article is as follows: An overview of previous studies on RAG models, the LangChain framework, multiple-search agent systems, and difficulties in UL settings is given in Section 2, "Related work". The system architecture and implementation techniques for a scalable RAG QA system utilizing LangChain are described in depth in Section 3, "Methodology," with an emphasis on the integration of multiple data sources through multi-search agents. The section titled "Evaluation" includes case studies that demonstrate the system's usefulness and practical applications. The system's flexibility, ability to grow, and impact on education and study are examined in Section 5, "Results and discussion." According to Section 6, "Conclusion," which also includes suggestions for future improvements in intelligent quality systems for educational tools, the results have been summed up. Additionally, it stresses the advantages of the RAG method with many search agents.

## 2. RELATED WORKS

In this section, we focus on related work that provides the foundation for the development of advanced question-answer (QA) systems in ubiquitous learning environments. The emphasis is on the synergistic fusion of Retrieval-Augmented Generation (RAG) models with the Langchain framework, leveraging agents and tools to enhance information retrieval and answer generation.

## 2.1 Retrieval-Augmented Generation (RAG) Models

Retrieval-Augmented Generation (RAG) models, introduced by [8], represent a hybrid approach that combines the strengths of retrieval-based and generative models. While the retrieval component gets relevant documents from a large corpus, the generative component makes replies that are appropriate for the situation [9]. Using this combination, RAG models can make automated Q&A systems more accurate and useful by basing answers on information from outside sources. Adding important facts from an outside knowledge base to the language model's internal representation of information is how RAG models help reduce problems caused by hallucinating false or misleading information [9]. Once the query and documents are retrieved, the generator, an advanced language model, uses them as context to make an answer. Furthermore, an attention method makes it easier for the generator to pick out the most important data from the retriever's output [10].

## 2.2 Langchain Framework

LangChain is a framework for developing applications powered by Large Language Models (LLMs) [10]. From creation to deployment, it makes the whole lifecycle of LLM applications easier. For easier application development, LangChain provides open-source building blocks, third-party tool interfaces, and templates [11]. Tools like LangSmith [12] The LangServe tool can help check, watch, and rate chain applications for ongoing improvement. [13] quickly converts applications into API services for deployment. For building context-aware processes, LangChain also links language models to outside data sources [14]. Due to being an open-source framework, it offers free access and community help, allowing both experts and beginners to use it.

For better learning, companies can make information retrieval and response generation faster, more accurate, and more useful to the situation by incorporating LangChain into QA systems in ubiquitous learning environments.

### 2.2.1 Tools

LangChain provides a suite of toolkits designed to streamline the development of LLM applications [15]. Wrapping around complicated processes, these toolkits make it easier for coders to add complex features without having to understand complicated code [16]. Along with helping to make context-aware workflows like RAG, these toolkits also make sure that apps can give correct and useful answers by using outside knowledge bases and reducing model hallucinations.

### 2.2.2 Agents

Using a language model's features to quickly decide a series of acts is what agents in the LangChain framework are all about. For chains, the order of actions is set and hardcoded. Agents, on the other hand, use a language model as a reasoning engine to choose which actions to perform and when. [17]. Being able to react to user queries more quickly and easily is made possible by this dynamic decision-making process. Depending on the context and content of the question, agents can use a variety of tools and data sources to make decisions in real time [18]. Therefore, using agents greatly enhances the overall efficiency and effectiveness of LLM applications, allowing for smarter and more relevant answers.

### 2.2.3 Vector embedding

The LangChain framework depends on vector embedding, which turns text data into dense, continuous models that make retrieval and generation better [19]. It is possible for the model to better understand and handle information because these embeddings record semantic relationships [20]. Entering a query turns it into a vector, which is then compared to stored document vectors to find the most appropriate ones. This method uses semantic connections [21], guaranteeing the retrieval of contextually relevant documents, which enhances the answer quality of the generative component.

## 2.3 Ubiquitous Learning Environments

Ubiquitous learning leverages pervasive computing technologies to create flexible and context-aware educational experiences [22]. Advanced technology supports this method, which lets learning happen naturally in a variety of settings and situations. Researchers are working on making systems that help people keep learning and change based on the needs of students in different settings. Integration of advanced question-and-answer systems, like the suggested multi-search agent RAG application, is in line with the goals of global learning because it gives students instant access to complete, accurate, and relevant information. Through real-time support and personalized learning pathways, this integration improves the learning experience and makes sure that educational tools are always available, wherever they are needed.

## 2.4 Challenges of RAG Models

The learning process can be greatly improved by combining RAG models with AI-based teaching tools. Examples include older programs like PAT2Math [23] and AI-based math tutor frameworks [24] for smart tutoring tools by laying the groundwork. While these systems did try to help, they were limited by their static question libraries and rule-based algorithms, which made it harder for them to come up with multiple, relevant answers. RAG models, while superior to traditional Q&A methods, face challenges in combining multiple data sources due to differences in format, quality, and exposure, as well as the strain on system resources for real-time processing. The accuracy of the retriever is critical, as retrieving too few relevant documents can degrade answer quality, necessitating constant monitoring and improvement. A multi-agent system can improve the accuracy and comprehensiveness of results by querying and integrating data from multiple sources using several search agents. This solves several difficulties. By increasing the adaptability and relevance of AI-based teaching aids, this strategy works well in educational contexts. In ubiquitous learning environments that value adaptability, ease of use, and individualized instruction, RAG models and multi-agent systems provide a substantial improvement in intelligent Q&A systems, providing more dependable and valuable results.

## 3. METHODOLOGY

This section offers a thorough overview of our system's technological architecture and execution. Our system's primary goal is to provide accurate and intelligent Question & Answer (Q&A) capabilities that are tailored to each user's unique needs. The Retrieval-Augmented Generation (RAG) pipeline is designed to include an autonomous layer, which is part of the system to build a responsive question answering system. We clarify the concept of agents, their role in using Large Language Models (LLMs) for automatic inference, and tool selection, which is the core functionality of our system.

To automate logical reasoning and select the best tools, the agent makes use of Language Models (LLMs).

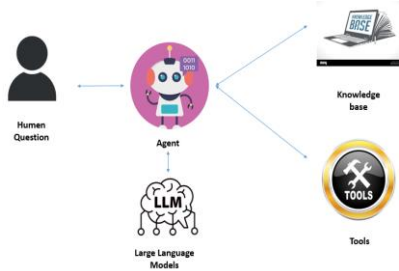


FIG. 1. THE OVERALL SOLUTION OF OUR SYSTEM

According to Figure 1, the main parts of the agent-based RAG system are:

- **Human Question:** The system initiates with a question posed by a user.
- **Agent:** The agent utilizes a Large Language Model (LLM) at the core of the system to perform automated reasoning and tool selection. The agent operates independently, making decisions and performing actions based on the specific circumstances of the interaction.
- **Knowledge Base:** The agent extracts pertinent information from an extensive knowledge base to guarantee precise and contextually suitable responses.
- **Tools:** The agent employs a range of tools to aid in the collection, analysis, and efficient presentation of information. This approach integrates the retrieval of information from the knowledge base with the generating capability of LLMs.

All of these parts work together to make an iterative, self-governing solution that provides users with complete and correct answers to their questions.

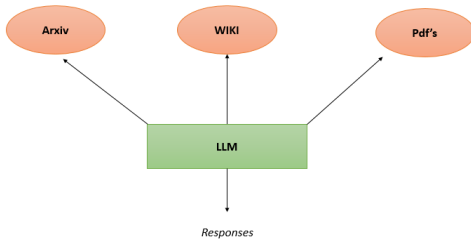


FIG. 2. WRAPPER METHOD

At its heart, the system depicted in Figure 2 is a Large Language Model (LLM) that communicates with three primary data sources: PDFs, Arxiv, and WIKI or Wikipedia. The graphic depicts the logic map's (LLM) processing of data received from various sources in order to provide a response. Multiple data sources can be easily integrated into the system by using a common wrapping technique [25], enabling their utilization; This system architecture relies on the wrapper's ability to ease access to several data sources. As a result, the LLM can do its job and give complete and accurate answers. Multiple sources provide us with our data.:

*Arxiv:* Refers to scholarly articles, usually sourced from the Arxiv repository, a renowned resource for research papers in disciplines such as physics, mathematics, computer science, and others.

*WIKI:* Refers to Wikipedia, an expansive internet-based compendium that offers extensive knowledge on a diverse array of subjects.

*PDFs:* Represents a diverse range of documents in PDF format, encompassing research papers, technical documents, reports, web pages, and other textual data.

## 2.1 System Architecture

The three primary parts of the proposed system design are the retrieval, generation, and integration modules. As seen in Figure 4, this design is tailored to optimally utilise and combine various data sources. An in-depth examination of each component:

### 3.1.1 Retrieval Module

Figure 3 shows the RAG workflow broken down into three separate steps: indexing, retrieval, and creation. Every phase's outcome has a major impact on the phase after it. Using the user's query (Q), the retrieval stage finds a subset of relevant documents, D (D1,..., Dk), from a collection of documents, E. P(D|Q) is the probability of choosing document D from set E in response to enquiry Q. It is the retriever's job to determine the probability. The next step, after obtaining the relevant documentation (D), is to create the solution (A). The probability distribution of producing the answer A, given the recovered documents D and the enquiry Q, can be expressed as P(A|D, Q), which represents the generation process. A text generator determines the probability of this happening.

The entire RAG process can be expressed as follows in formula (1):

$$P(A | Q) = \sum_{i=1}^K P(A | D_i, Q) \cdot P(D_i | Q) \quad (1)$$

The notation P(A|Q) represents the likelihood of producing answer A when presented with enquiry Q.

In conclusion, the RAG algorithm takes into consideration both the input context and the chosen documents in order to generate the output text. It is usual practice to employ models based on neural networks to do such operations. In order to maximise the probability of proper output and alignment with the retrieved texts, these models are trained. When comparing two vectors A and B, the cosine similarity is employed as a measure of their similarity. The following is the formula (2) for it:

$$\text{Cosine Similarity}(A, B) = \frac{A \cdot B}{\|A\| \cdot \|B\|} = \frac{\sum_{i=1}^n A_i \cdot B_i}{\sqrt{\sum_{i=1}^n A_i^2} \cdot \sqrt{\sum_{i=1}^n B_i^2}} \quad (2)$$

The three separate steps that make up the RAG workflow are indexing, retrieval, and generation (Figure 4). Each phase's outcome greatly affects the next consecutive step:

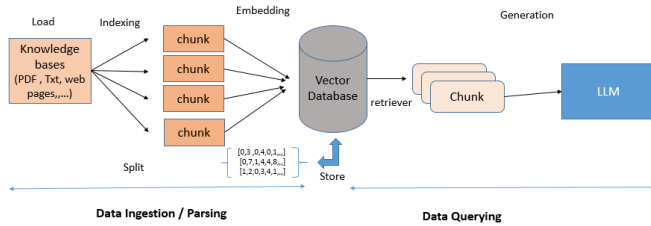


FIG. 3. PROTOTYPING NAIVE RAG TOOL.

**Indexing:** We are indexing the data to make retrieval as efficient as possible. The process of indexing helps to arrange information in a way that allows for quick searching and retrieval.

**Splitting:** A manageable portion of the indexed data is separated into smaller parts. The success of the data processing and handling of large documents depends on this stage.

**Embedding:** A vector representation is thereafter used to incorporate each segment. Embeddings convert chunks into numerical formats that capture their semantic meaning, which is essential for similarity-based retrieval.

**Vector Database:** The embedded chunks are stored in a vector database. This specialized database is designed to handle vectorized data, enabling fast and accurate retrieval based on semantic similarity.

### 3.1.2 Generation Module

The Generating module employs a pre-trained language model, such as GPT-4, that has been specifically adjusted for Q&A jobs. This model produces responses by utilising the documents obtained by the retrieval module, guaranteeing that the answers are both precise and contextually appropriate. The generation module use Language Models (LLMs) to generate responses that are contextually precise and logically connected. After the retrieval module obtains pertinent documents, they are inputted into the LLM, which utilises the information to produce a thorough response. The LLM is utilized by the agent to initially process the user's query and ascertain the most suitable line of action.

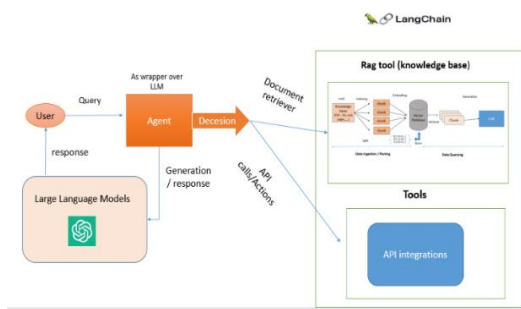


FIG. 4. THE OVERALL ARCHITECTURE OF OUR SYSTEM

### 3.1.3 Integration Module

The integration module merges the results from the retrieval and generating modules, refining and synthesizing the information to generate a conclusive solution. This module is responsible for managing user interactions and feedback, and it enhances the system's performance by making iterative modifications.

- **Agent:** The agent is the central node and building block of the system because of its role as the primary decision-maker. Automated reasoning and tool selection both make advantage of the LLM. The agent can draw on its knowledge base, come up with an answer independently, or use external tools via APIs.
- **Tools (API Integrations):** The agent can access a variety of tools using API integrations provided by the integration module. Based on the user's request, these tools assist in carrying out targeted actions, collecting supplementary data, or gaining access to up-to-the-minute information.

As shown in Figure 4, the process begins with the user interacting with the system by asking a query. Upon receiving the enquiry, the agent promptly analyses it to determine the most appropriate response. The decision-making process involves weighing the benefits of directly constructing a response, connecting to other technologies through APIs, or collecting pertinent data from the database. Once the agent chooses the right approach, all of the collected and generated data is combined to provide a final answer. Sending this back to the user completes the interaction cycle. This design guarantees a versatile, self-sufficient, and iterative solution that provides users with complete and accurate answers by efficiently combining retrieval-augmented generation with multi-agent systems.

### 2.2 Implementation

In this section, our methodology offers a systematic and comprehensive way to construct a state-of-the-art RAG system, guaranteeing the incorporation of various data sources, honing language models, and creating a unified framework for precise and pertinent question and answer-generation.

Developing the retrieval-augmented generation (RAG) pipeline requires multiple phases using Python 3.10.0 and LangChain:

- **Data Source Integration:** Ensuring a varied range of information requires incorporating multiple data sources into the retrieval module. Make use of wrapper tools for integrating datasets, such as Wikipedia, and follow Figure 5's instructions. Aside from that, we will display custom data using web sources. Chunking, embedding with OpenAI embeddings, and storing in a vector store (FAISS) are all part of the process. We went with FAISS because of its similarity search capabilities in sets of vectors of any size and its support for local index storage. By avoiding its creation at each iteration and instead loading it directly, we may reduce processing time. Figure 6 depicts the code for the retriever, and Figure 7 shows the creation of the ArXiv algorithm. You can see all three of the RAG system's tools in Figure 8.

```
from langchain_community.tools import WikipediaQueryRun
from langchain_community.utilities import WikipediaAPIWrapper
api_wrapper=WikipediaAPIWrapper(top_k_results=1,doc_content_chars_max=200)
wiki=WikipediaQueryRun(api_wrapper=api_wrapper)
```

FIG.5. CREATE A WIKIPEDIA TOOL



FIG.6. CREATE A WEBSITE RETRIEVER TOOL

FIG.7. CREATE THE ARXIV TOOL

FIG.8. THE TOTAL OF TOOLS USED

- FIG.9. LLM MODEL USED

FIG.10. CREATE PROMPT

- FIG.11. CREATE AGENT

FIG.12. AGENT EXECUTOR

In this section, we describe how the system processes a user query to generate an accurate and relevant response. When a user submits a query, the agent first analyzes it to determine the appropriate course of action. This may involve searching a vector database for relevant documents, retrieving full documents from a document store, calling external APIs, performing computations, or employing a combination of these actions. As information is gathered, the agent refines its internal representation of the query. The LLM then generates a response based on the refined query and the retrieved information, which is subsequently sent back to the user.

Several case studies in ubiquitous learning scenarios highlight the practical benefits of the system. To illustrate how various tools and approaches are used to handle user requests, we have included some examples below.

Figure 13 displays the results of our study, which show how we used the Retriever Augmented Generation (RAG) approach to get useful information out of a vector database.

FIG.13. THE RESULT OF CASE STUDY 1

1. *Query Submission*: The user enters a question, such as "Please tell me more about LangSmith."
2. *Document Retrieval*: The agent searches the vector database for papers about LangSmith using the Retriever tool.
3. *Information Processing*: The papers that have been retrieved are processed to extract the relevant parts.
4. *Response Generation*: The LLM then uses the updated query and the data it has retrieved to construct a response.
5. *User Response*: The resulting response is delivered back to the user, which includes relevant case studies for LangSmith.

In this case, the study presented in Figure 14, showcases how the system uses the Wikipedia API to fetch general information related to a user query.

FIG. 14. THE RESULT OF CASE STUDY 2

6. *Query Submission*: The user queries, "Explain the concept of ubiquitous learning."
  7. *API Call*: The agent calls the Wikipedia API to search for articles related to ubiquitous learning.
  8. *Document Retrieval*: The Wikipedia API returns a relevant article on ubiquitous learning.
  9. *Information Processing*: The agent processes the article to extract key points and a concise explanation of the concept.
  10. *Response Generation*: The LLM generates a detailed response explaining ubiquitous learning based on the extracted information.
  11. *User Response*: The explanation is sent back to the user, providing them with a clear understanding of the concept of ubiquitous learning.
- 2.3.3 *Case Study 3: Using Arxiv for Insights into Deep Learning Techniques*

In this case, the study presented in Figure 15, we illustrate how the system leverages the Arxiv API to provide comprehensive insights into the latest advancements in deep learning techniques.

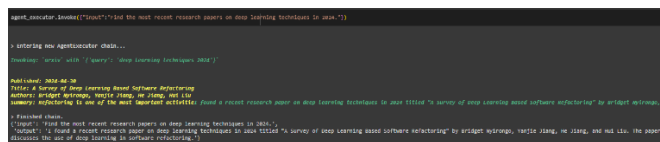


FIG. 15. THE RESULT OF CASE STUDY 3

1. *Query Submission*: The user submits a query, "Find the most recent research papers on deep learning techniques in 2024."
2. *API Call*: The agent utilizes the Arxiv API to search for the most recent research papers published on the topic of deep learning techniques.
3. *Document Retrieval*: The Arxiv API returns a list of the most relevant and recent research papers matching the query criteria.
4. *Information Processing*: The agent processes the abstracts, key sections, and conclusions of the retrieved papers to extract important information and insights.
5. *Response Generation*: The LLM generates a response summarizing the key findings, innovations, and trends in deep learning techniques based on the extracted information.
6. *User Response*: The summary, including insights from the latest research papers on deep learning techniques, is sent back to the user.

## 5. RESULTS AND DISCUSSION

In this section, the performance of our system is thoroughly examined in this presentation, with numerous case studies serving as examples. The system's flexibility, efficiency, and adaptability across various applications are highlighted when we go over the wider consequences of these results. Read on for an outline of the most important takeaways from implementing our system, which demonstrates its versatility, speed, and capacity to process massive datasets while also facilitating ongoing research and learning. By showcasing the improved performance and user happiness acquired through the integration of various tools

and APIs within a coherent framework, the system's potential in academic research, professional development, and general education is demonstrated through its successful implementation across various scenarios. Below is a summary of table 1 of the results from each case study.

The outcomes showcase the system's remarkable adaptability and versatility. Academic research, professional development, and general education are just a few areas that can benefit from the system's versatility and efficiency in answering a wide variety of questions. The agent can handle a broad variety of user needs, providing both high-level overviews and in-depth analyses, thanks to its usage of several methodologies and techniques. Quick data collection and processing from large databases and APIs is evidence of the system's efficiency and scalability. To efficiently manage ever-increasing data volumes and ensure that users can swiftly get the most relevant and up-to-date information, this capability is crucial. Users in search of comprehensive information find the system invaluable since it uses vector databases and retrievers like FAISS to improve its power to search and retrieve pertinent documents.

TABLE I Summary of Case Study Results

Case Study	Tool/Method Used	Query	Outcome	Implications
Case Study 1	Retriever (RAG)	"tell me about LangSmith."	Detailed information was retrieved and summarized from a vector database.	Demonstrates efficient handling of specific queries, valuable for users seeking detailed information from large datasets
Case Study 2	Wikipedia API	"Explain the concept of ubiquitous learning."	Clear and concise explanation of ubiquitous learning.	Effectiveness in providing accurate information for educational and informational purposes.
Case Study 3	Arxiv API	"Deep learning techniques."	Recent research papers on deep learning techniques fetched and summarized.	Aids professionals and researchers in staying informed about the latest developments, supporting continuous learning and innovation.

The system facilitates research and learning by providing quick access to the latest research findings and comprehensive summaries, helping researchers and professionals stay current and fostering innovation. Its ability to efficiently gather and analyze complex research from platforms like Arxiv highlights its potential for supporting advanced research and continuous knowledge acquisition. Experiments demonstrate that the multi-search agent RAG system excels in delivering accurate and contextually appropriate answers by utilizing diverse data sources, thereby overcoming the limitations of traditional Q&A systems. This approach enhances the system's resilience and reduces biases associated with single-source retrieval.

**Improved Answer Quality and Contextuality:** The contextual relevance agent checks that the retrieved data is

suitable for the given context, leading to more relevant and high-quality results.

The decision-making agent chooses the best data sources and tools on the go, so the system can change to suit different kinds of enquiries and user requirements. The overall performance of the system and user happiness are both improved by this flexibility.

The system's capacity to adapt and handle multiple query patterns and data sources is demonstrated by its successful implementation in these distinct case studies. Dynamic and contextually relevant information can be retrieved and generated by integrating numerous tools and APIs into a single Framework. This capability has significant implications for fields like as education, research, and knowledge dissemination, where accurate and up-to-date information is crucial.

## 6. CONCLUSION

The study introduces a Retrieval-Augmented Generation (RAG) system that is both dynamic and agent-based. It draws on a variety of data sources to provide users with thorough and accurate responses to their queries. Automated reasoning and tool selection are powered by Large Language Models (LLMs), making the system versatile enough to be used in general education, professional development, and academic research. To guarantee the transmission of high-quality, contextually relevant information, its architecture is designed to incorporate retrieval, generation, and integration components. Improved system performance and user happiness are the results of the system's multi-agent framework's ability to dynamically choose suitable tools and data sources. The system's efficiency, scalability, and adaptability are showcased in case studies, particularly when it comes to rapidly obtaining relevant information from big databases. Dynamic adaptation by the decision-making agent enhances system efficacy, while multiple data sources are integrated to assure unbiased replies.

Some potential improvements for the future include adding more tools and APIs, moving to a cloud-based architecture for greater scalability, making decisions better based on user feedback, adding support for new languages, creating easier interfaces, and improving the breadth of tools and APIs. We also suggest advanced security techniques to safeguard user data. All things considered, the RAG system is a huge step forward in the realms of query responding and data retrieval, and it can completely change the way people in all kinds of industries deal with data.

## REFERENCES

- [1] M. Guettala, S. Bouekkache, and O. Kazar, "Ubiquitous learning a new challenge of ubiquitous computing: state of the art," in *2021 International Conference on Information Systems and Advanced Technologies (ICISAT)*, 2021, pp. 1-5.
- [2] M. Guettala, S. Harous, S. Bouekkache, B. Athamena, O. Kazar, and Z. Houhamdi, "Cloud ubiquitous learning approach based on multi-agents system," in *2022 International Arab Conference on Information Technology (ACIT)*, 2022, pp. 1-8.
- [3] X. Zhao, "Research on Methods and Applications Related to Question-and-Answer Dialogue Systems," *Highlights in Science, Engineering and Technology*, vol. 57, pp. 9-14, 2023.
- [4] S. Basu, A. S. Rawat, and M. Zaheer, "A statistical perspective on retrieval-based models," in *International Conference on Machine Learning*, 2023, pp. 1852-1886.
- [5] D. M. Anstine and O. Isayev, "Generative models as an emerging paradigm in the chemical sciences," *Journal of the American Chemical Society*, vol. 145, pp. 8736-8750, 2023.
- [6] J. Chen, H. Lin, X. Han, and L. Sun, "Benchmarking large language models in retrieval-augmented generation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024, pp. 17754-17762.
- [7] S.-Q. Yan, J.-C. Gu, Y. Zhu, and Z.-H. Ling, "Corrective retrieval augmented generation," *arXiv preprint arXiv:2401.15884*, 2024.
- [8] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, and T. Rocktäschel, "Retrieval-augmented generation for knowledge-intensive nlp tasks," *Advances in Neural Information Processing Systems*, vol. 33, pp. 9459-9474, 2020.
- [9] H. Li, Y. Su, D. Cai, Y. Wang, and L. Liu, "A survey on retrieval-augmented text generation," *arXiv preprint arXiv:2202.01110*, 2022.
- [10] W. Yu, D. Iyer, S. Wang, Y. Xu, M. Ju, S. Sanyal, C. Zhu, M. Zeng, and M. Jiang, "Generate rather than retrieve: Large language models are strong context generators," *arXiv preprint arXiv:2209.10063*, 2022.
- [11] R. Asyofi, M. R. Dewi, M. I. Lutfhi, and P. Wibowo, "Systematic Literature Review Langchain Proposed," in *2023 International Electronics Symposium (IES)*, 2023, pp. 533-537.
- [12] T. Ito, T. Kuribayashi, M. Hidaka, J. Suzuki, and K. Inui, "Langsmith: An interactive academic text revision system," *arXiv preprint arXiv:2010.04332*, 2020.
- [13] D. Gao, Z. Li, W. Kuang, X. Pan, D. Chen, Z. Ma, B. Qian, L. Yao, L. Zhu, and C. Cheng, "AgentScope: A Flexible yet Robust Multi-Agent Platform," *arXiv preprint arXiv:2402.14034*, 2024.
- [14] A. K. Jadoon, C. Yu, and Y. Shi, "ContextMate: a context-aware smart agent for efficient data analysis," *CCF Transactions on Pervasive Computing and Interaction*, pp. 1-29, 2024.
- [15] J. Jin, Y. Zhu, X. Yang, C. Zhang, and Z. Dou, "FlashRAG: A Modular Toolkit for Efficient Retrieval-Augmented Generation Research," *arXiv preprint arXiv:2405.13576*, 2024.
- [16] W. Sun, J. Wang, Q. Guo, Z. Li, W. Wang, and R. Hai, "CEBench: A Benchmarking Toolkit for the Cost-Effectiveness of LLM Pipelines," *arXiv preprint arXiv:2407.12797*, 2024.
- [17] M. Abbasian, I. Azimi, A. M. Rahmani, and R. Jain, "Conversational health agents: A personalized llm-powered agent framework," *arXiv preprint arXiv:2310.02374*, 2023.
- [18] P. Neira-Maldonado, D. Quisi-Peralta, J. Salgado-Guerrero, J. Murillo-Valarezo, T. Cárdenas-Arichábal, J. Galan-Mena, and D. Pulla-Sanchez, "Intelligent Educational Agent for Education Support Using Long Language Models Through Langchain," in *International Conference on Information Technology & Systems*, 2024, pp. 258-268.
- [19] X. Huan and H. Zhou, "Integrating Advanced Language Models and Vector Database for Enhanced AI Query Retrieval in Web Development," *International Journal of Advanced Computer Science & Applications*, vol. 15, 2024.
- [20] Q. Liu, M. J. Kusner, and P. Blunsom, "A survey on contextual embeddings," *arXiv preprint arXiv:2003.07278*, 2020.
- [21] J. Kandola, N. Cristianini, and J. Shawe-taylor, "Learning semantic similarity," *Advances in Neural Information Processing Systems*, vol. 15, 2002.
- [22] M. Guettala, S. Bouekkache, O. Kazar, S. Harous, and M. Zouai, "The design and implementation of intelligent ubiquitous learning Multi-Agent Context-Aware System," *World Journal on Educational Technology: Current Issues*, vol. 15(4), pp. 429-450, 2023.
- [23] P. A. Jaques, H. Seffrin, G. Rubi, F. de Moraes, C. Ghilardi, I. I. Bittencourt, and S. Isotani, "Rule-based expert systems to support step-by-step guidance in algebraic problem solving: The case of the tutor PAT2Math," *Expert Systems with applications*, vol. 40, pp. 5456-5465, 2013.
- [24] W. Gan, Y. Sun, S. Ye, Y. Fan, and Y. Sun, "AI-tutor: Generating tailored remedial questions and answers based on cognitive diagnostic assessment," in *2019 6th international conference on behavioral, economic and socio-cultural computing (BESC)*, 2019, pp. 1-6.