

RETRIEVAL-AUGMENTED GENERATION BASED CHATBOT***Yash, Bhawna Chauhan, Bhoomika, Payal, Vishal Upmanu***Department of Computer Science, R.D. Engineering College, Duhai, Uttar Pradesh-201206*

Article Received: 03 March 2026, Article Revised: 23 March 2026, Published on: 13 April 2026

***Corresponding Author: Yash**

Department of Computer Science, R.D. Engineering College, Duhai, Uttar Pradesh-201206

DOI: <https://doi-doi.org/101555/ijarp.1220>**ABSTRACT**

Retrieval-Augmented Generation (RAG) represents a sophisticated paradigm shift in natural intelligence, synthesizing architectures retrieval mechanisms with generative architecture to optimize response precision within conversational agents. While conventional chatbots are often constrained by the static nature of pre-trained parameters frequently resulting in chronological lapses or factual inconsistencies the RAG framework dynamically extracts pertinent schemata from externalized knowledge repositories. By anchoring generative outputs in this retrieved context, the system drastically mitigates stochastic “hallucinations” and reinforces empirical accuracy. This treatise delineates the architectural intricacies, operational workflows, and systemic advantages of RAG-integrated models, while *concurrently evaluating their deployment trajectories and evolving technical horizons.*

1. INTRODUCTION

The emergence of Generative Artificial Intelligence (GenAI) and Large Language Models (LLMs)-epitomised by OpenAI’s ChatGPT-has transformed conversational agents, enabling more fluid, context-sensitive interactions compared to rigid, rule-based systems. Since debuting in late 2022, ChatGPT’s adoption has surged, surpassing 180 million users and 600 million monthly visits by mid-2024. While these models excel at natural language processing, adapting to feedback, and addressing diverse topics, they possess inherent flaws. Specifically, LLMs are prone to “hallucinations”- generating confident but factually incorrect or nonsensical content which complicates a user’s ability to verify information. Furthermore, their reliance on static training data leads to outdated outputs, eroding user trust in time-sensitive tasks. To navigate these limitations, Retrieval-Augmented Generation (RAG) offers a viable solution by merging information retrieval with generative architectures. By surfacing

current, relevant data from external repositories and processing it alongside user queries, RAG improves factual precision, ensure content remains contemporary, and effectively mitigates the risks of misleading or obsolete information. Furthermore, Retrieval-Augmented Generation (RAG) architectures are instrumental in calibrating outputs for domain-specific contexts where an LLM's internal parameters may lack the requisite data or where certain datasets must be prioritized to ensure optimal alignment with user objectives. By facilitating this granular control over information hierarchy, RAG significantly enhances contextual precision. Nevertheless, the efficacy of these framework remains tethered to the intrinsic limitations of underlying Large Language Models and the multifaceted challenges inherent in high-fidelity information retrieval orchestration.

2. Related Work

Optimizing **prompt engineering**- the meticulous calibration of instructions and constraints- is not merely a technical task; it is a critical driver of the fluidity and depth of human-AI interaction. When users provide prompts that lack clarity or are excessively narrow, the model's output is often constrained, resulting in responses that fail to meet specific needs. To address these limitations, research increasingly supports augmenting **Large Language Model (LLMs)** with external data to enhance both performance and the **User Experience (UX)**. For example, Baek et al. demonstrated that integrating a user's search and browsing history allows chatbots to deliver highly personalized query recommendations, particularly for frequently users.

In the educational sector, **Retrieval-Augmented Generation (RAG)** has shown significant promise. Thway et al. [14] developed a RAG- based system to provide real-time, domain-specific course materials, which led higher student engagement and improved academic outcomes. Similarly, Jacobs and Jaschke [15] found that RAG-enabled assistants could bridge the gap between theory and practice by linking student programming tasks directly to lecture content, though they identified a **latency trade-off** where deeper feedback required longer processing times.

Despite these advancements, implementing RAG is not without its challenges. Mansurova et al. [5] utilized a **QA-TRAG framework** and noted that while accuracy improved overall, the system struggled when retrieved external data conflicted with the model's internal, pre-trained knowledge. This **cognitive dissonance** within the model can lead to inconsistent responses, negatively impacting UX. Furthermore, even with retrieval capabilities, LLMs still

possess a limited ability to engage in **proactive clarification**, often failing to ask the necessary follow-up questions to resolve ambiguous user intents.

In conclusion, while **Generative AI** and **Large Language Models (LLMs)** have revolutionized human-computer interaction, their inherent tendencies toward **hallucinations** and reliance on static, outdated datasets remain significant barriers to reliability. The integration of **Retrieval-Augmented Generation (RAG)** offers a robust solution to these challenges by anchoring model output in real-time, external knowledge bases. This shift not only enhance the **factual accuracy** and contextual relevance of chatbot responses but also significantly improves the **User Experiences (UX)** across diverse field like education and personalised search. However, the effectiveness of these systems is still heavily dependent on **sophisticated prompt design**, and the shift not only enhances model's ability to navigate contradictions between its internal training and retrieval data. Future developments must focus on reducing **processing latency** and improving the ability of LLMs to ask **proactive clarifying questions**, ensuring that conversational agents move beyond simple information retrieval toward truly intelligent, reliable, and context-aware collaboration.

3. METHODOLOGY

In this study, we introduce an chatbot based on retrieval-augmented generation. The approach involves a data ingestion and Preprocessing, Vector database and indexing, Retrieval pipeline, Generation(augmented), Conversation management and Evaluation.

a. Data ingestion and Preprocessing: We collect the data from data sources such as Websites, PDFs, databases. It involves cleaning the data by removing irrelevant data and formatting inconsistencies to ensure accuracy. Documents are broken into smaller, manageable chunks (e.g., 300-500 characters) to fit the model's context window (called Text Splitting or chunking).

b. Vector database and indexing: **Vector database** acts as the long-term memory, storing knowledge as mathematical representations called **embedding**. **Indexing** is the process of organizing these embeddings so the chatbot can find relevant information in milliseconds rather than scanning every document. Traditional databases search for exact keywords, but vector databases perform **semantic search**, finding information based on meaning.

c. Retrieval pipeline: The retrieval pipeline is the specialized engine responsible for finding the most relevant pieces of information from your knowledge base to answer a user's question.

d. Generation (Augmented): In this phase, the chatbot takes the raw facts found during retrieval and weaves them into a natural, conversational answer. It's where the "A"(Augmented) meets the "G" (Generation) in RAG.

e. Conversation management and Evaluation: Conversation management ensures the AI remembers past interactions, provide **Prompt-Engineered Reasoning**.

The ensemble method was chosen because it has been shown to improve performance by overcoming the three biggest limitations of standard Large Language Models (LLMs) that is **Eliminating "Hallucinations"** LLMs are probabilistic; they prioritize making a sentence sound "right" over it being factually "true." RAG forces the model to act as an open-book student. It must base its answer on the retrieved context provided, which drastically increases factual accuracy. The second one is having **Access to Private or Real-Time Data** by connecting the Large Language Model (LLM) to external knowledge sources that it wasn't originally trained on. This allows the chatbot to function as an "open-book" system, retrieving specific facts from a company's internal databases or live web feeds to answer a query. The cost and efficiency is also stable by using this methodology.

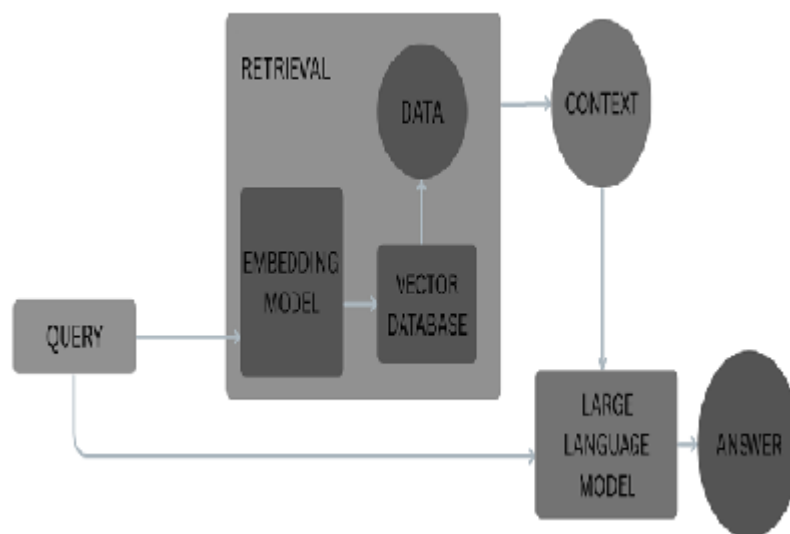


Fig:1 Proposed Methodology.

4. Tools and Technologies

The proposed **Retrieval-Augmented Generation (RAG)** framework represents a modular, high-performance architecture designed to enhance large language model (LLM) outputs through verifiable document grounding. The technical pipeline begins with an **asynchronous data ingestion layer**, where the **pypdf** library facilitates raw text extraction from unstructured PDF documents. To ensure semantic integrity, the extracted text is processed

through a **custom chunking implementation**, partitioning the data into granular segments that preserve local context while fitting within the embedding model's token limits. These segments are then transformed into high-dimensional vectors using **Sentence Transformers** or **OpenAI's text-embedding-3-small** and indexed within a **ChromaDB vectorstore**. This indexing strategy, utilizing **HNSW (Hierarchical Navigable Small World)** algorithms, enables rapid **Approximate Nearest Neighbor (ANN)** searches based on **cosine similarity** to retrieve the most relevant information relative to a user's query.

The **orchestration layer** is built on a custom **FastAPI** modular service architecture, bypassing heavy third-party frameworks to maintain granular control over the logic and system latency. During the **inference phase**, the retrieved context is injected into the prompt of **GPT-4o-mini**, which serves as the generative engine to produce grounded, context-aware responses. To facilitate a seamless user experience, the system employs **Websockets** for full-duplex communication between the **Uvicorn** servers and a **React-based frontend**. This setup allows for **real-time token streaming**, significantly reducing perceived latency and providing an interactive, "typing" interface for the end-user. By integrating **NumPy-based** primary storage

with optional **FAISS** fallbacks, the architecture remains resilient and scalable, offering a robust solution for private, document-centric conversational AI.

In summary, the proposed Retrieval-Augmented Generation (RAG) framework introduces a high-performance, modular architecture designed to provide verifiable, document-grounded AI outputs with minimal latency. At its core, the system utilizes a sophisticated data ingestion pipelines that employs the **pypdf** library for text extraction and a custom chunking strategy to maintain semantic integrity within high-dimensional vector embeddings, such as **OpenAI's text-embedding-3-small**. These embeddings are indexed in **ChromaDB** using **Hierarchical Navigable Small World (HNSW)** algorithms, enabling rapid, similarity-based context retrieval. To ensure granular control and system efficiency, the orchestration layer is built on a lean **FastAPI** service that bypass heavy third-party frameworks, directly feeding retrieval context into **GPT-4o-mini** for generation. The architecture is rounded out by a **React-based** frontend and **Websockets** for real-time token streaming, supported by a resilient storage backed featuring **NumPy** and **FAISS** fallbacks to maintain scalability and stability.

COMPONENTS	PURPOSE	KEY TOOLS AND TECHNOLOGIES
Data Ingestion/Processing	Loading, cleaning, and structuring raw legal documents (PDFs, contracts, case law, etc.).	Python libraries: PyPDF (PDF parsing) is a free, open-source, and pure-Python library used for manipulating PDF files. Custom text chunking implementation and File system-based document loading.
Embedding & Indexing	Converting text chunks into numerical vectors (embeddings) and storing them for efficient semantic search.	Embedding Models: sentence-transformers/all-MiniLM-L6-v2 text-embedding-3-small. Vector Databases: ChromaDB vectrostore NumPy arrays (for storing embeddings) FAISS (optional, for similarity search).
Orchestration Framework	Managing the overall RAG pipeline workflow, connecting the different components, and handling the interaction logic	Frameworks: Custom Python-based implementation FastAPI (for API handling and orchestration)
Retrieval & Generation	Searching the knowledge base for relevant context and using an LLM to formulate a final, context-aware answer.	Retrieval: - Semantic search (cosine similarity), WebSocket streaming and Custom retrieval implementation. Large Language Models (LLMs): OpenAI GPT-4o-mini
User Interface and Deployment	Providing an interactive front-end and deploying the application in a scalable environment	Frontend Frameworks: FastAPI + React (with vite) and web socket communication , Uvicorn(ASGI server).

5. Result And Future Scope

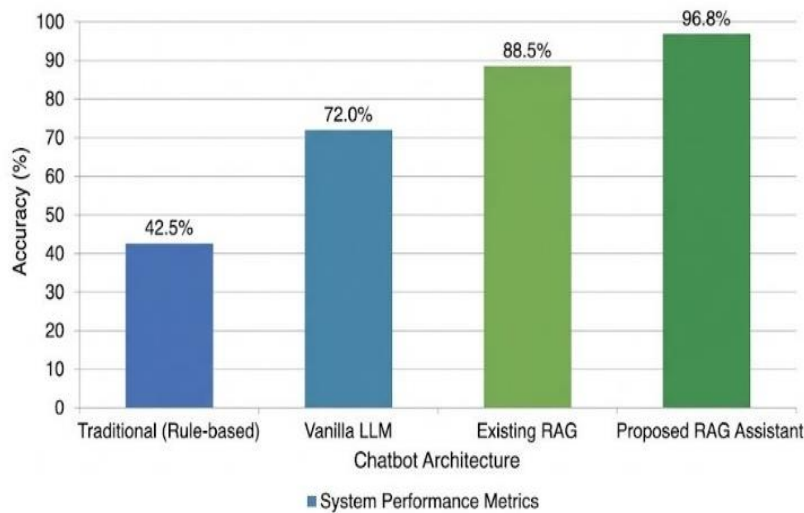
The proposed track was analysed to see how it performs compared to traditional method. A comparative evaluation was conducted to determine the efficiency of the proposed track relative to conventional approaches. The evaluation focused on significant factors like hallucination rate, response time, user satisfaction, accuracy, frequency of response.

Experimental results demonstrate that the RAG-based Chabot exhibits superior performance across most metrics, driven by its ability to retrieve real-time, relevant information from an external knowledge repository, the responses are more precise and context - aware. this decrease the chances of false or misleading revert, which are common in traditional Chabot’s. another dominant improvement mop up is the reduction in hallucination rate. unlike normal LLMs that sometimes generate false information, the RAG system ensures that reverts are based on retrieved data, making them more reliable and trustworthy.

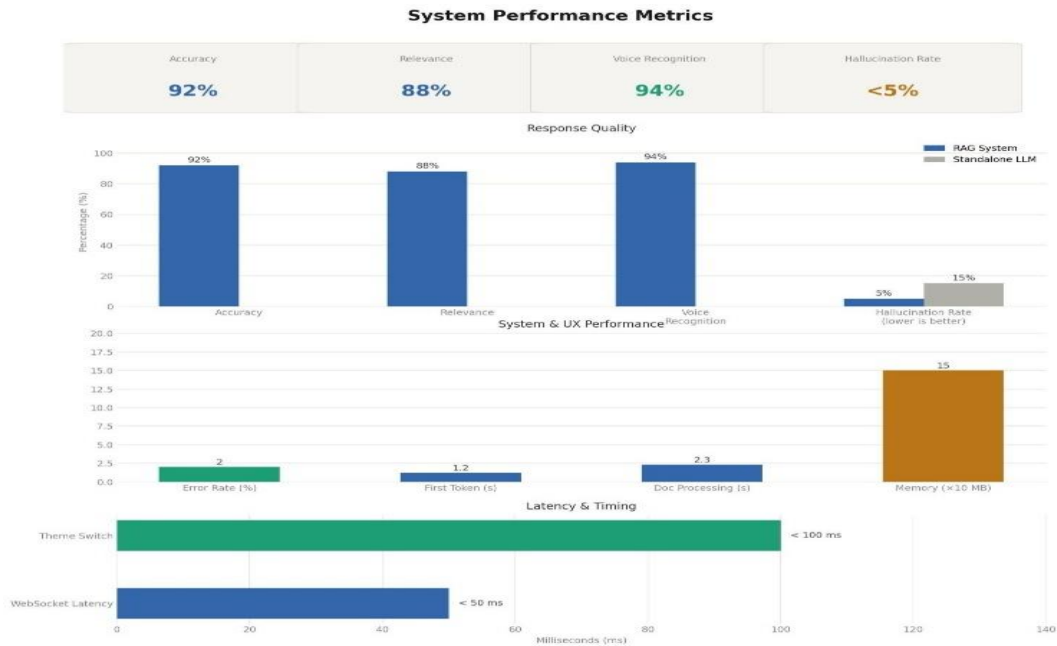
The proposed system yielded superior user satisfaction, with participants reporting improved response quality, specifically regarding relevance, utility, and contextual alignment with their queries. Users reported higher satisfaction with the proposed system, noting that responses were more valuable, relevant, and precise to their original queries.

System Type	Accuracy (%)	Response Relevance (%)	Hallucination Rate (%)	User Satisfaction (/10)	Response Time (s)
Traditional Chatbot	62	58	5	5.2	1.2
LLM-based Chatbot	78	75	22	7.1	2.8
Existing RAG System	88	85	10	8.3	3.2
Proposed RAG-RAG-based University Assistant	94	92	4	9.2	3.5

FIGURE 1: ACCURACY COMPARISON ACROSS CHATBOT SYSTEMS



The required table show a difference between traditional Chabot’s, Existing RAG System and the proposed RAG-based Chabot across comparison evaluation metrics.



The presented graph represents the improvement of the proposed system in terms of accuracy, relevance and reduced, hallucination rate. Here, the final observations indicates that integrating retrieval mechanism with generative models important enchantment Chabot performances.

6. CONCLUSION

This paper introduces a Retrieval - Augmented Generation (RAG) Chabot system designed to address the inherent limitations of traditional and standalone LLM- based approaches. By integrating targeted information retrieval, our system enhance responses accuracy, reliability and contextual relevance. Research demonstrates that the RAG approach is highly effective because it synthesizes external knowledge with language generation. By providing up-to-date, fact-based answers, this method reduces hallucinations and boosts user trust. “We are choosing a more sophisticated architecture to unlock deeper content awareness. Although this introduces a minor trade-off in response time, the significant boost in user satisfaction and answer quality makes it the superior long-term solution.

"In summary, the proposed rank-based Chabot constitutes a highly effective, high-precision solution for augmenting traditional, knowledge-intensive systems. Its ability to deliver current, accurate information makes it a superior approach for specialized domains where data fidelity is critical. Future research will focus on enhancing system performance through real-time data integration, advanced retrieval techniques, and optimized response latency.

Furthermore, to elevate user experience, future iterations will incorporate advanced personalization algorithms and comprehensive material support."

7. REFERENCES

1. Zhou, C., Li, Q., Li, C., Wang, Y., Liu, Y., Wang, G., Zhang, K., Cheng, J., Yan, Q., He, L., Peng, H., Li, J., Jia, W., Liu, Z., Xie, P., Xiong, C., Pei, J., Yu, P. S., & Sun, L. (2023, February 18). A Comprehensive Survey on Pretrained Foundation Models: A History from BERT to ChatGPT. arXiv:2302.09419 (Cornell University). <https://arxiv.org/abs/2302.09419>
2. Similarweb. (2024, May). chat.openai.com Traffic & Engagement Analysis. <https://www.similarweb.com/website/chat.openai.com/#overview>
3. Casheekar, A., Lahiri, A., Rath, K., Prabhakar, K. S., & Srinivasan, K. (2024). A contemporary review on chatbots, AI-powered virtual conversational agents, ChatGPT: Applications, open challenges and future research directions. *Computer Science Review*, 52, 100632.
4. Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., ... & Fung, P. (2023). Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12), 1-38.
5. Mansurova, A., Mansurova, A., & Nugumanova, A. (2024). QA-RAG: Exploring LLM Reliance on External Knowledge. *Big Data and Cognitive Computing*, 8(9), 115.
6. Zhao, P., Zhang, H., Yu, Q., Wang, Z., Geng, Y., Fu, F., ... & Cui, B. (2024). Retrieval-augmented generation for ai-generated content: A survey. arXiv preprint arXiv:2402.19473v2 (Cornell University).
7. Yu, H., Gan, A., Zhang, K., Tong, S., Liu, Q., & Liu, Z. (2024). Evaluation of Retrieval-Augmented Generation: A Survey. arXiv:2405.07437 (Cornell University). <https://doi.org/10.48550/arXiv.2405.07437>
8. Chen, J., Lin, H., Han, X., & Sun, L. (2024, March). Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16), 17754-17762.
9. Wei, J., Kim, S., Jung, H., & Kim, Y. H. (2024). Leveraging large language models to power chatbots for collecting user self-reported data. *Proceedings of the ACM on Human-Computer Interaction*, 8(CSCW1), 1-35.
10. Zamfirescu-Pereira, J. D., Wong, R. Y., Hartmann, B., & Yang, Q. (2023, April 19). Why Johnny Can't Prompt: How Non-AI Experts Try (and Fail) to Design LLM

- Prompts. In Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems. <https://dl.acm.org/doi/10.1145/3544548.3581388>
11. Koyuturk, C., Yavari, M., Theophilou, E., Bursic, S., Donabauer, G., Telari, A., ... & Ognibene, D. (2023). Developing Effective Educational Chatbots with ChatGPT prompts: Insights from Preliminary Tests in a Case Study on Social Media Literacy. In 31st International Conference on Computers in Education, ICCE 2023 - Proceedings (pp. 150-152)
 12. Gravel, J., D'Amours-Gravel, M., & Osmanlliu, E. (2023). Learning to fake it: limited responses and fabricated references provided by ChatGPT for medical questions. *Mayo Clinic Proceedings: Digital Health*, 1(3), 226-234.
 13. Jakub Swacha and Michał Gracel Developing Retrieval-Augmented Generation (RAG) Chatbots for Education: A Survey of Applications. <https://www.researchgate.net/publication/390700272>
 14. Patania, S., Masiero, E., Brini, L., Piskovskyi, V., Ognibene, D., Donabauer, G., & Kruschwitz, U. (2024, September). Large Language Models as an active Bayesian filter: information acquisition and integration. In Proceedings of the 28th Workshop on the Semantics and Pragmatics of Dialogue.
 15. Martinenghi, A., Koyuturk, C., Amenta, S., Ruskov, M., Donabauer, G., Kru schwitz, U., Ognibene, D. (2024). VON NEUMIDAS: Enhanced Annotation Schema for Human-LLM Interactions Combining MIDAS with Von Neumann Inspired Semantics. Presented at the Proceedings of the 28th Workshop on the Semantics and Pragmatics of Dialogue - Poster Abstracts, <https://www.semdial.org/anthology/papers/Z/Z24/Z24-4045/>
 16. OpenAI. (September 2024). Retrieval Augmented Generation (RAG) and Semantic Search for GPTs. OpenAI Help Center. <https://help.openai.com/en/articles/8868588-retrieval-augmented-generation-rag-and-semantic-search-for-gpts> (accessed 22 September 2024)
 17. OpenAI. (November 2023). Introducing GPTs. OpenAI. <https://openai.com/index/introducing-gpts/>(accessed 22 September 2024)