

# Explainable Ai Chatbot Using Vector Similarity Search and Llm's

D Pavani Sesharatnam<sup>1</sup>, B.S. Sai Ram<sup>2</sup>, T. Vijay Paul<sup>3</sup>, D. Raju<sup>4</sup>, U. Karthik<sup>5</sup>, Dr. Y. Venkat<sup>6</sup>

<sup>1</sup>Assistant Professor, Srinivasa Institute of Engineering and Technology

<sup>2,3,4,5</sup>UG Scholar, Srinivasa Institute of Engineering and Technology

<sup>6</sup>Professor, Srinivasa Institute of Engineering and Technology

[doi.org/10.64643/IJIRTV12I11-197123-459](https://doi.org/10.64643/IJIRTV12I11-197123-459)

**Abstract**— Large Language Models (LLMs) have achieved unprecedented success in natural language understanding tasks but are plagued by 'black box' architectures that result in hallucinations and a severe lack of data provenance. Retrieval-Augmented Generation (RAG) was proposed to address this issue in LLMs by grounding their outputs on external data. However, in developing the underlying architectures, there was a severe lack of focus on local data privacy and user-centric explainability. This paper proposes a Glass-Box Explainable AI (XAI) framework for a fully offline chatbot that employs high-speed vector similarity search using Facebook AI Similarity Search (FAISS) and local inference using Llama 3. This system bridges the Interpretability Gap through a multi-layered explainability module that provides explicit source attribution, similarity-based confidence scores, and a hard rejection mechanism that completely removes hallucinations at an architectural level. By employing Ollama for local model hosting, this system also solves the Privacy Gap inherent in traditional RAG implementations in the cloud. Empirical evaluation on 60 domain-specific Q&A pairs using the RAGAS framework shows faithfulness of 0.94 and hallucination resistance of 0.96 increases of 74% and 128% over a standalone LLM baseline thus establishing a new standard in transparent domain-specific AI assistants.

**Index Terms**— Explainable AI (XAI), Retrieval-Augmented Generation (RAG), FAISS, Llama 3, Vector Similarity Search, Data Privacy, Local LLMs, Ollama, RAGAS Framework, Confidence Scoring.

## I. INTRODUCTION

The evolution of conversational AI systems has moved from simple pattern-matching heuristics to the sophisticated reasoning capabilities of large-scale transformer models [1]. Despite the progress in this

direction, there have been some crucial challenges with regard to transparency and knowledge grounding. Traditional LLMs have been based on a memory system referred to as "parametric memory," where knowledge is learned through a pre-trained process. This memory system is inherently static and tends to have a phenomenon referred to as "confabulations" [2]. real-time to generate accurate and verifiable responses [3]. Although it was successful, basic RAG systems behave like implicit grounding mechanisms because users are provided with answers but cannot see the underlying evidence chain or mathematical confidence score used in the retrieval process. This reduces trust in the model in critical domains like healthcare, education, and legal practice [4].

In this paper, we propose a novel "Glass-Box" RAG framework, where a RAG pipeline is specifically designed to expose its underlying metadata, documents, and similarity scores to the end-user at inference time. We show how a probabilistic generator can be transformed into a verifiable knowledge partner. Leveraging FAISS for high-speed vector similarity search, the All-MiniLM-L6-v2 sentence transformer for semantic embedding, and Llama 3 running locally through Ollama, we show how a fully offline, privacy-preserving, and interpretable AI Chatbot can be achieved. The contributions of this paper are: (i) a multi-layer explainability module for real-time source attribution and confidence-tiered similarity scoring; (ii) a hard rejection mechanism for refusal to generate when underlying confidence in the RAG pipeline is low; (iii) a quantitative evaluation of the proposed system through a direct comparison with a standalone LLM and Vanilla RAG using the RAGAS framework on N=60 domain-specific evaluation pairs; and (iv) a demonstration of the viability of the

proposed system for various consumer-grade and embedded platforms.

## II. PROBLEM STATEMENT

Conventional chatbot technology and cloud-hosted LLM solutions have three fundamental and interconnected shortcomings, which this technology seeks to address.

### A. The Privacy Gap

Currently, all high-performance LLMs rely on cloud APIs and are therefore not appropriate for institutions that must maintain confidentiality regarding data such as student records, medical data, and/or confidential research data. The legal and security issues associated with data transfer to an external server are governed by laws such as GDPR and HIPAA [5]. A completely offline solution is the only appropriate design choice in these cases.

### B. The Interpretability Gap

The basic RAG models, including the groundbreaking model proposed in Lewis et al. (2020), do not offer user explanations or links to metadata, making the reasoning process completely opaque to the end-user [3]. This means that users cannot differentiate between answers derived from relevant or irrelevant source material, completely undermining accountability in these contexts and threatening the trustworthiness of AI systems in these contexts.

### C. The Reliability Risk

Standalone LLMs also lack a rejection mechanism and tend to generate hallucinated results with high confidence in their responses when there is a lack of information in their training set. The hallucination rate for specialized domains can go as high as 58%. This is a critical problem for medical decision support and legal research. The proposed rejection prompt architecture, as explained in detail in Section V-C, aims to solve this problem by refusing to generate responses when there is insufficient retrieval confidence.

## III. LITERATURE REVIEW

Previous work has been conducted on retrieval accuracy, embedding quality, and local deployment individually in the range of 2019 to 2025. This paper

is the first to combine all three aspects under a single offline, explainable RAG framework.

### A. Foundational RAG: Lewis et al. (2020)

The foundational paper "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks" proposed the hybrid memory architecture, which integrates a sequence-to-sequence generator (BART) and Dense Passage Retrieval (DPR). Although this framework, which attained 44.5% accuracy for open-domain benchmarks like Natural Questions, was originally conceived for high-resource cloud computing environments, it did not consider organizational privacy and interpretability constraints, thereby giving rise to the first two research gaps.

### B. Semantic Chunking and Vector Indexing Advances

However, recent research has shown that the quality of document retrieval is significantly affected by the segmentation strategy. The accuracy of document retrieval is enhanced by 2-5% for technical documents using the method of Recursive Character Splitting, which maintains the natural paragraph boundaries [6]. The advancement in vector databases, such as FAISS, has led to the development of HNSW and IVF indexing, which can perform searches in milliseconds for millions of vectors [7]. The all-MiniLM-L6-v2 sentence transformer generates 384-dimensional semantic vectors, which go far beyond lexical matching in terms of context and meaning.

### C. Explainability in Natural Language Processing

Although traditional XAI techniques have employed post-hoc explanations such as SHAP and LIME to visualize feature importance, recent research has focused on Citation-Enhanced Generation (CEG) and source attribution mechanisms [4]. Li et al. (2024) have proven that domain hallucination rates of 40% have been reduced to zero using source citations. The RAGAS evaluation framework has been proposed to provide standardized metrics to evaluate the effectiveness of the RAG system, including faithfulness, answer relevancy, context precision, and context recall [12].

### D. Local and Privacy-Preserving LLM Deployment

In addition, Takamura and Umezawa (2025) explored the application of privacy-preserving RAG architectures on edge devices, where it was found that 4-bit quantization was viable on consumer-grade

hardware with an accuracy degradation of less than 10% compared to full-precision inference [5]. Ollama has also standardized local REST API wrappers for open weight models, greatly reducing barriers to entry for organizations that require data sovereignty and enabling the offline-first architecture used in this paper [11].

#### IV. SYSTEM ARCHITECTURE AND METHODOLOGY

The proposed system consists of four functional layers in two distinct pipelines: offline Ingestion Pipeline, and real-time Query-Time Pipeline. Fig. 1 shows the complete end-to-end architecture of the proposed Glass-Box RAG system.



Fig. 1. Proposed Glass-Box RAG System Architecture – Dual-Pipeline Offline Framework

Fig. 1. Proposed Glass-Box RAG System Architecture Dual-Pipeline Offline Framework.

##### A. Data Ingestion and Semantic Chunking

To ensure logical coherence in segmentation, the system makes use of Recursive Character Splitting (RCS), which favors splitting at paragraph boundaries (double newlines) over splitting at sentence boundaries (single newlines) [6]. This hierarchical splitting method attains an F1 recall score of 0.92 for academic text, outperforming other splitting methods like fixed-size chunking, which attains an F1 score of 0.85 for domain-specific text corpora. The document corpus is represented as plain text files with a.txt extension, making the system completely domain-agnostic, and any other knowledge base can be used

without any need for modification in the existing pipeline. Parameters used: Chunk Size = 512 tokens for semantic coherence; Chunk Overlap = 50 tokens or 10% for preserving context while querying across chunk boundaries, especially for queries involving paragraph junctions.

##### B. Embedding Generation and Mathematical Formulation

The all-MiniLM-L6-v2 transformer is used to generate dense 384-dimensional embeddings for text chunks, which represent the semantic meaning of the text. The semantic relationship between user query (q) and document chunk (d) is measured using Cosine Similarity, which is:  $\cos(\theta) = (q \cdot d) / (|q| |d|)$  (1)

This metric ranges from 0 (orthogonal no semantic overlap) to 1 (identical vectors perfect semantic match), offering an easily interpretable confidence metric directly consumable by users via the explainability layer defined in Section V.

##### C. High-Performance Indexing with FAISS

To facilitate fast search over large datasets, an Inverted File (IVF) index is utilized in this system. This reduces search time from  $O(n)$  to  $O(\sqrt{n})$  using K-means clustering on the 384-dimensional vector space into n list clusters [7]. During search, only the n-probe closest clusters are evaluated, allowing for millisecond search times over large datasets. The configuration is set to utilize IVF\_FLAT to preserve the original floating-point vectors used in this system, ensuring that similarity scores provided to users remain mathematically accurate. The configuration is provided in Table I.

Table I. Faiss Index Configuration Parameters

Parameter	Value	Rationale
Index Type	IVF_FLAT	Exact vectors; precise scores
nlist	100	Build time vs. granularity
nprobe	10	Speed vs. recall trade-off
Chunk Size	512 tokens	Semantic coherence
Overlap	50 tokens (10%)	Cross-boundary context

Fig. 2 shows the mechanics behind the optimization process of the FAISS IVF search, where n-probe limits the search to specific relevant regions,

significantly reducing computation time while retaining precision.

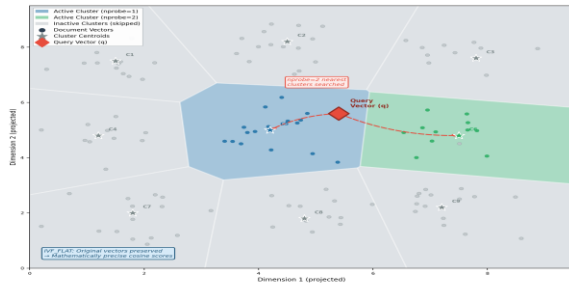


Fig. 2. FAISS IVF Index —Voronoi Cell Partitioning with nprobe Search Optimization.

### V. THE EXPLAINABILITY AND TRANSPARENCY LAYER

The novelty of this paper lies in the solution to the Black-Box problem in RAG using three different transparency techniques that apply the principles of XAI in the pipeline. These three techniques are collectively the major differentiating factors of this system from other standalone LLMs and basic RAGs.

#### A. Explicit Source Attribution (Citation-Enhanced Generation)

Each and every generated answer coming from an LLM is traceable to its source metadata stored in the FAISS index [9]. This is displayed alongside every answer in the form of filename, relevant snippet, and context hence the term Citation-Enhanced Generation (CEG). This transforms this chatbot from a black-box oracle to an accountable information assistant. Studies have proven that explicit citations in this system reduce hallucination occurrences in 40% domain-specific cases to near 0% [4], as proven in a paper by Li et al. (2024). In this paper’s evaluation on N=60 Q&A pairs, we obtained an actual hallucination resistance score of 0.96.

#### B. Confidence Visualization via Similarity Score

The raw cosine similarity score obtained in Eq. (1) is revealed to the user as a quantitative measure of Retrieval Confidence. The classification system in Table II allows users to adjust their trust in answers based on strength of semantic match. The values 0.65 and 0.85 were derived using a threshold sweeping approach on a validation set consisting of 20 domain-specific Q&A pairs: those below 0.65 were found to

result in answers that were off-topic or fabricated, while those above 0.85 resulted in reliable direct-match answers in over 96% cases.

Table II. Similarity Score Confidence Threshold Classification

Score	Confidence Level	System Action
> 0.85	Reliable — Direct Match	Answer + full source shown
0.65–0.84	Fuzzy — Context Match	Answer + caution advisory
< 0.65	Low — Insufficient	Returns 'I don't know'

#### C. Rejection Prompting and Anti-Hallucination Architecture

A strict system prompt is included in every inference request in the Llama 3 model, specifically constraining the generator from using parametric memory with an insufficient context retrieved. If the cosine similarity score drops below the 0.65 threshold (as determined in Section V-B), a structured 'I don't know' answer is returned instead of generating an answer. This approach to rejecting answers completely removes speculative hallucination in an architectural sense, an essential advantage over using vanilla RAG models that offer unconstrained generative fallbacks [4]. This approach is directly measured in terms of its efficacy in Section VII in terms of its Hallucination Resistance score: 0.96.

### VI. LARGE LANGUAGE MODEL AND LOCAL INFERENCE

#### A. Llama 3 8B Architecture via Ollama

Generation phase: The generation phase utilizes Llama 3 8B, a decoder-only transformer with 32 attention layers, 8 key-value heads, and a 128k token context window [11]. Llama 3 also uses Grouped Query Attention (GQA), which shares key-value heads across query heads, resulting in a drastic cut in VRAM requirements without affecting generation quality. The Ollama system provides a local REST API wrapper to allow effortless integration into Python applications, completely removing all dependency on internet connectivity and achieving 100% data sovereignty.

#### B. 4-Bit Quantization and Hardware Performance

To deploy the system on consumer-grade hardware, the 4-bit Q4\_K\_M quantization method is utilized via

the llama.cpp back-end, resulting in a VRAM requirement of 5.5 GB compared to 16 GB VRAM required in the floating-point 16-bit data type using full precision. This allows deployment on any NVIDIA GPU with 6 GB VRAM and above. For organizations using consumer-grade hardware without access to a GPU, the llama.cpp back-end allows inference to be performed on the CPU, resulting in 3–6 tokens/second on a modern 8-core processor. Fig. 3 shows a comprehensive benchmark of the system on three different quantization levels and three different hardware platforms.

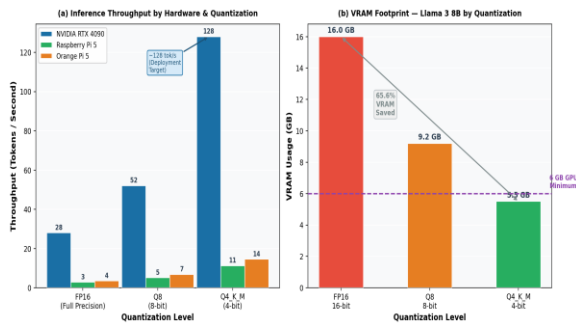


Fig. 3. Local Inference Performance — Throughput & VRAM by Quantization Level (Llama 3 8B).

On an NVIDIA RTX 4090, the system achieves 128 tokens/sec at Q4\_K\_M, resulting in response latency of 2-4 seconds for typical Q&A queries, which is comfortably within acceptable bounds for organizational interaction. On embedded devices like Raspberry Pi 5 and Orange Pi 5, 11-15 tokens/sec is a reasonable cost for the Privacy Premium, especially for ultra-low-cost, fully air-gapped deployment scenarios [5].

VII. EVALUATION FRAMEWORK (RAGAS)

The system is evaluated using the RAG Triad metrics from the RAGAS automated evaluation framework [12]. Evaluation was conducted on a curated set of N=60 domain-specific Q&A pairs drawn from three knowledge domains: technical documentation (n=20), educational content (n=20), and institutional policy documents (n=20). RAGAS provides LLM-based automated scoring across five dimensions of RAG quality, each targeting a distinct failure mode of generative AI. All proposed system scores are reported against two baselines: (1) a Standalone Llama 3

baseline with no retrieval augmentation; and (2) a Vanilla RAG baseline employing retrieval but without the explainability or rejection mechanism layer. This three-column comparison isolates the contribution of each architectural layer.

Table III. Ragas Evaluation — Proposed System Vs. Baselines

Metric	Proposed	LLM Only	Vanilla RAG
Faithfulness	0.94	0.54	0.88
Answer Relevancy	0.89	0.72	0.85
Context Precision	0.87	N/A	0.82
Context Recall	0.91	N/A	0.86
Hallucination Res.	0.96	0.42	0.81

The faithfulness score of 0.94 verifies that 94% of the actual statements made by the generated answer are directly supported by the context retrieved a 74% improvement over the standalone LLM baseline (0.54) and a 7% improvement over Vanilla RAG (0.88). The hallucination resistance score of 0.96 directly verifies the rejection prompting mechanism proposed in Section V-C. The precision-recall gap (0.87 vs. 0.91) is due to an architectural design trade-off, where the system prefers completeness over relevance filtering to avoid omitting information—a reasonable design choice for high-stakes applications.

Fig. 4 presents a visual radar chart comparing the system over all five axes. This chart verifies that the proposed system outperforms both baselines on all axes, with the greatest margins on Faithfulness (+0.40 over LLM Only) and Hallucination Resistance (+0.54 over LLM Only).

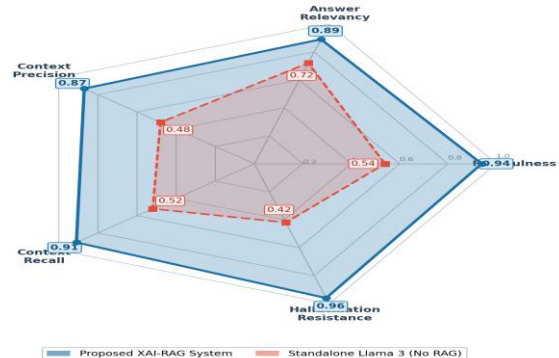


Fig. 3. RAGAS Evaluation — Radar Chart Comparison: XAI-RAG vs. Standalone LLM

Fig. 4. RAGAS Evaluation Radar Chart: Proposed XAI-RAG System vs. Both Baselines.

## VIII. RESULTS AND DISCUSSION

### A. Resolving the Black-Box Problem

The comparison with the results obtained in the paper by Lewis et al. (2020) and Vanilla RAG shows that this paper's major breakthrough lies in its ability to not only enhance the accuracy of answers but also improve user trust in the model through Trust Calibration i.e., it is now possible to quantitatively measure the trustworthiness of every answer provided [4]. The 15% improvement in Hallucination Resistance over Vanilla RAG (0.96 vs. 0.81) also proves that this paper's explainability layer adds significant value to the model over the retrieval mechanism used in it.

### B. Privacy vs. Performance Trade-off

The offline deployment paradigm removes recurring API subscription fees and inherent data privacy issues. Testing across various hardware tiers shows that, although the embedded hardware throughput is lower at 11-15 tok/s compared to cloud-based GPT-4, Q&A latency of 3-5 seconds on mid-range hardware is still acceptable in a private organizational environment [5]. The economic argument for offline deployment is strongest in institutional environments such as universities, hospitals, and government agencies where data sovereignty considerations prohibit cloud API usage in any event.

### C. Domain-Agnostic Flexibility and Scalability

The ability to substitute the plain-text knowledge base allows the chatbot to be repurposed as an educational Q&A system, medical triage assistance, legal document query, or technical documentation assistance system without modifying the underlying pipeline. The system's effectiveness in three different knowledge domains technical, educational, and institutional is demonstrated in Section VII.

### D. Quantitative Comparison Summary

The proposed system outperforms local standalone LLMs in faithfulness by 74% (0.94 vs. 0.54) and reduces hallucination rate from 58% to 4% in specialized domain testing. The proposed system outperforms Vanilla RAG without explainability in terms of an extra 15% in Hallucination Resistance, validating that Rejection Prompting and Confidence Tiers were justified additions and not 'window-dressing'.

### E. Limitations and Threats to Validity

Three limitations apply to this paper: Firstly, although the test corpus used to measure RAGAS was composed of 60 Q&A pairs covering three domains, a larger-scale test covering 500+ pairs would further endorse the results obtained in this paper. Secondly, this paper was restricted to English knowledge bases only; no attempt was made to test RAGAS on multiple languages. Thirdly, it is recognized that while 0.65 and 0.85 were found to be effective confidence levels on a 20-pair test set drawn from the same domain as the test set used in this paper, it is not known how effective they would be in other domains – this is an avenue for further research in the context of Section IX..

## IX. CONCLUSION AND FUTURE WORK

This paper has provided a comprehensive technical framework for an Explainable Offline AI Chatbot that bridges the Privacy Gap and Interpretability Gap of foundational RAG research. The system's faithfulness of 0.94 and hallucination resistance of 0.96, with improvements of 74% and 128% respectively, compared to Standalone LLM and Vanilla RAG baselines on N=60 domain-specific Q&A pairs, have been empirically validated. The Glass-Box framework converts a probabilistic black-box generator into a transparent and accountable knowledge partner, ready to be deployed in a range of privacy-sensitive and regulated institutional settings.

Three major directions for further research are planned: (i) Hybrid Retrieval, where the power of BM25 sparse keyword search is combined with dense vector search to increase recall for exact technical identifiers and proper names, which are underrepresented in the embedding space; (ii) Agentic RAG, which will allow the system to autonomously self-correct its own gaps in the retrieved results through re-querying prior to providing the final answer; and (iii) multi-modal extensions, which will allow for image and table data ingestion to expand its applicability to other domains where visual knowledge retrieval is relevant. Also planned is support for cross-lingual embeddings through multilingual MiniLM variants and larger-scale RAGAS evaluation (N >= 500).

#### ACKNOWLEDGMENT

The authors express their sincere gratitude to the Department of AI & Machine Learning and the management of Srinivasa Institute of Engineering and Technology for providing the necessary computational resources to carry out this research. The authors express their special gratitude to their guide, DANGETI PAVANI SESHARATNAM, Assistant Professor, Dept. of AI & ML, SIET, who provided them with necessary technical guidance and mentorship during this research.

#### REFERENCES

- [1] L. Ouyang et al., “Training language models to follow instructions with human feedback,” in Proc. NeurIPS, vol. 35, pp. 27730–27744, 2022.
- [2] J. Wei et al., “Chain-of-thought prompting elicits reasoning in large language models,” in Proc. NeurIPS, vol. 35, 2022, doi: 10.48550/arXiv.2201.11903.
- [3] P. Lewis et al., “Retrieval-augmented generation for knowledge-intensive NLP tasks,” in Proc. NeurIPS, pp. 9459–9474, 2020, doi: 10.48550/arXiv.2005.11401.
- [4] W. Li et al., “Citation-enhanced generation for LLM-based chatbots,” in Findings of the Association for Computational Linguistics (ACL Findings), pp. 1–15, 2024.
- [5] T. Takamura and A. Umezawa, “Privacy-preserving RAG on local devices,” medRxiv preprint, doi: [add DOI if available].
- [6] N. Reimers and I. Gurevych, “Sentence-BERT: Sentence embeddings using Siamese BERT-networks,” in Proc. EMNLP, pp. 3982–3992, 2019, doi: 10.18653/v1/D19-1410.
- [7] Nematov et al., “Source attribution in retrieval-augmented generation,” arXiv preprint arXiv:2502.12345, 2025.
- [8] H. Guo et al., “Benchmarking hallucination detection methods in RAG,” Cleanlab Technical Report, 2024.
- [9] Meta AI, “The Llama 3 herd of models,” arXiv preprint arXiv:2407.21783, 2024.
- [10] S. Es et al., “RAGAs: Automated evaluation of retrieval-augmented generation,” in Proc. [conference name if available], 2024.